

# Near-Optimal Stochastic Bin-Packing in Large Service Systems with Time-Varying Item Sizes

YIGE HONG, Carnegie Mellon University, USA

QIAOMIN XIE, University of Wisconsin-Madison, USA

WEINA WANG, Carnegie Mellon University, USA

In modern computing systems, jobs' resource requirements often vary over time. Accounting for this temporal variability during job scheduling is essential for meeting performance goals. However, theoretical understanding on how to schedule jobs with time-varying resource requirements is limited. Motivated by this gap, we propose a *new setting* of the stochastic bin-packing problem in service systems that allows for *time-varying* job resource requirements, also referred to as 'item sizes' in traditional bin-packing terms. In this setting, a job or 'item' must be dispatched to a server or 'bin' upon arrival. Its resource requirement may vary over time while in service, following a Markovian assumption. Once the job's service is complete, it departs from the system. Our goal is to minimize the expected number of active servers, or 'non-empty bins', in steady state.

Under our problem formulation, we develop a job dispatch policy, named JOIN-REQUESTING-SERVER (JRS). Broadly, JRS lets each server independently evaluate its current job configuration and decide whether to accept additional jobs, balancing the competing objectives of maximizing throughput and minimizing the risk of resource capacity overruns. The JRS dispatcher then utilizes these individual evaluations to decide which server to dispatch each arriving job to. The theoretical performance guarantee of JRS is in the asymptotic regime where the job arrival rate scales large linearly with respect to a scaling factor  $r$ . We show that JRS achieves an additive optimality gap of  $O(\sqrt{r})$  in the objective value, where the optimal objective value is  $\Theta(r)$ . When specialized to constant job resource requirements, our result improves upon the state-of-the-art  $o(r)$  optimality gap. Our technical approach highlights a novel policy conversion framework that reduces the policy design problem into a single-server problem.

CCS Concepts: • **Networks** → **Network performance analysis**; • **Mathematics of computing** → **Markov processes**.

Additional Key Words and Phrases: stochastic bin-packing, large service systems, policy conversion, asymptotic optimality

## ACM Reference Format:

Yige Hong, Qiaomin Xie, and Weina Wang. 2023. Near-Optimal Stochastic Bin-Packing in Large Service Systems with Time-Varying Item Sizes. *Proc. ACM Meas. Anal. Comput. Syst.* 7, 3, Article 48 (December 2023), 46 pages. <https://doi.org/10.1145/3626779>

## 1 INTRODUCTION

### 1.1 Background and Motivation

In modern computing systems, a job often takes the form of a virtual machine (VM) or a container [8, 14]. Such a job comes with a resource requirement, such as a certain number of CPUs and amount of memory, while in service. Each server in the system offers a limited amount of these resources.

Authors' addresses: Yige Hong, [yigeh@andrew.cmu.edu](mailto:yigeh@andrew.cmu.edu), Carnegie Mellon University, Pittsburgh, Pennsylvania, USA, 15213; Qiaomin Xie, [qiaomin.xie@wisc.edu](mailto:qiaomin.xie@wisc.edu), University of Wisconsin-Madison, Madison, Wisconsin, USA, 53706; Weina Wang, [weinaw@cs.cmu.edu](mailto:weinaw@cs.cmu.edu), Carnegie Mellon University, Pittsburgh, Pennsylvania, USA, 15213.



This work is licensed under a Creative Commons Attribution International 4.0 License.

© 2023 Copyright held by the owner/author(s).

2476-1249/2023/12-ART48

<https://doi.org/10.1145/3626779>

When a job arrives at the system, the job dispatch policy needs to decide which server the job should be assigned to, given the job's resource requirement and servers' current job configurations. This job scheduling problem can be approached as a *Stochastic Bin-Packing (SBP) problem*, where jobs are viewed as items, job resource requirements as item sizes, and servers as bins. A traditional SBP setting considers a finite set of jobs that arrive online but do not depart from the system. The objective is to minimize the number of servers that have jobs on them, or 'non-empty bins', subject to the resource capacities of the servers. SBP, with a rich history in operations research and theoretical computer science [10–12], is a field of continuous developments and advancements [1, 16, 18].

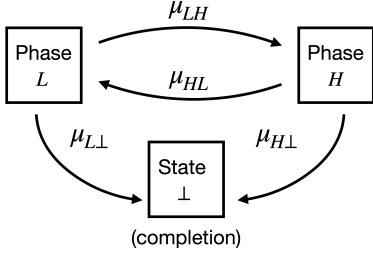
To incorporate job *departures* into the problem formulation, a setting referred to as *stochastic bin-packing in service systems* has been proposed recently [17, 32–36]. In this setting, jobs not only arrive but also depart over time. More specifically, jobs are assumed to arrive according to Poisson processes, and each job is assumed to stay in the system for an exponentially distributed service time. The service time of a job remains unknown until the job departs. Before delving further into SBP in service systems, it is worth mentioning that there is a parallel thread of research on the so-called dynamic bin-packing problem that also handles job departures (see, e.g., [6, 9, 20], and references therein), but it is primarily from a worst-case analysis perspective. Additionally, the virtual machine scheduling problem with objectives different from minimizing the number of active servers has also been widely studied (see, e.g., [22–24, 26–29, 39]).

For SBP in service systems, the goal is to design a job dispatch policy  $\sigma$  that minimizes the expected number of active servers in *steady state*, denoted as  $N(\sigma)$ . The *optimality gap* of a policy  $\sigma$  is defined as  $N(\sigma) - N(\sigma^*)$ , where  $\sigma^*$  is the optimal policy. Since SBP in service systems aims to model today's large-scale computing systems, the optimality gap of a policy is usually studied in the regime where the total job arrival rate becomes large. As we scale up the total job arrival rate linearly with a scaling factor,  $r$ , the optimal value  $N(\sigma^*)$  can be shown to be  $\Theta(r)$ .<sup>1</sup> Therefore, we say a policy is asymptotically optimal if its optimality gap is  $o(r)$ .

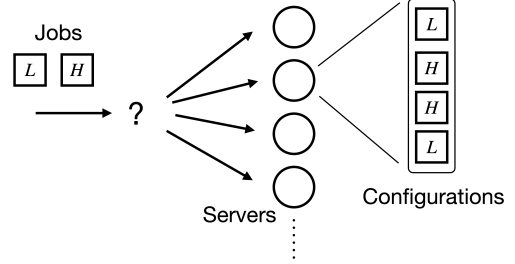
The optimality gap for SBP in service systems has been characterized in the line of work [32–36]. In particular, Stolyar [32], Stolyar and Zhong [34] propose greedy policies that are asymptotically optimal, but the scheduler that executes these policies needs to know detailed state information, which is in a high-dimensional space. Stolyar and Zhong [35, 36] later develop policies that use much less state information and achieve  $\Theta(r)$  (with an arbitrarily small constant) and  $o(r)$  optimality gaps, respectively.

While prior work on SBP in service systems has provided substantial insights into scheduling virtual-machine-type jobs, it primarily focuses on job resource requirements that remain fixed over time. However, in modern computing systems, jobs' resource requirements often vary over time [2, 13, 21, 30, 31, 37]. For example, when a job involves providing user-facing services, the instantaneous requirement on CPUs and memory depends on the service demand, which is subject to fluctuation over time [13, 21]. Time-varying job resource requirements pose significant challenges in optimizing system efficiency, particularly when aiming to minimize the number of active servers, thereby improving server utilization. It is pertinent to note that low utilization has been recognized as a significant obstacle to the continued scaling of today's computing systems.

Motivated by this gap, in this paper, we propose a *new setting of stochastic bin-packing in service systems* that allows job resource requirements, or 'item sizes', to vary over time.



(a) A simplified version of our job model. Each job in service is in either an  $L$  phase or an  $H$  phase, associated with low and high resource requirements, respectively. When the job is completed, it is said to be in the state  $\perp$ . The job transitions between the two phases while in service until it is completed, following a continuous-time Markov chain with rates  $\mu_{ii'}$ ,  $i, i' \in \{L, H, \perp\}$ .



(b) A system model with an infinite number of identical servers. As soon as a job arrives to the system, the job needs to be dispatched to a server to start service immediately. The configuration of each server is the number of jobs in each phase on the server.

Fig. 1. Job model and system model.

## 1.2 Problem Formulation: A Simplified Version

We first describe our job model that features time-varying resource requirements. For ease of exposition, here we present a simplified setting where each job in service can be in one of the two *phases*,  $L$  and  $H$ , associated with *low* and *high* resource requirements, respectively. Our full model, presented in Section 2, allows *more than two phases and more than one type of resources*. To model the temporal variation in the resource requirement, we assume that each job transitions between the two phases while in service until it is completed, following a continuous-time Markov chain illustrated in Figure 1(a). We use an absorbing state  $\perp$  to denote that the job is completed. A job can initialize in either phase  $L$  or phase  $H$ , and they are referred to as *type L* and *type H* jobs, respectively. Note that the setting where a job's resource requirement does not vary over time is a special case of our job model where the transition rates between phases are 0.

We consider a system with an infinite number of identical servers, illustrated in Figure 1(b). We assume jobs arrive according to a Poisson process as existing work on SBP in service systems. In particular, we assume that the two types of jobs arrive at the system following two independent Poisson processes, with rates  $\Lambda_L$  and  $\Lambda_H$ , respectively; i.e., the interarrival times of type  $L$  and type  $H$  jobs are i.i.d. following exponential distributions with means  $1/\Lambda_L$  and  $1/\Lambda_H$ , respectively. Upon arrival, a job needs to be dispatched to a server according to a *dispatch policy*, and the job enters service immediately. The goal is to design a policy  $\sigma$  to minimize the expected number of *active servers* (servers currently serving a positive number of jobs) in steady state, denoted as  $N(\sigma)$ .

As job resource requirements vary over time, situations can arise where the total job resource requirement on a server exceeds the server's resource capacity, resulting in resource contention. Modern computing systems can tolerate temporary overruns of resource capacity, though they often incur performance degradation or other costs [7, 15]. In our model, we incorporate a rate at which the cost accumulates due to resource contention. We first represent the state of a server by its *configuration*, a vector  $\mathbf{k} = (k_L, k_H)$  where  $k_L$  and  $k_H$  are the numbers of jobs in phase  $L$  and phase  $H$ , respectively. Then a *cost rate function*  $h(\cdot)$  maps a server's configuration to a rate of cost. For example, the cost rate can be proportional to how much the total resource requirement of

<sup>1</sup>We use the standard Bachmann–Landau notation. Consider two functions  $a(r)$  and  $b(r)$ , where  $b(r)$  is positive for large enough  $r$ . Then  $a = O(b)$  if  $\limsup_{r \rightarrow +\infty} \frac{|a(r)|}{b(r)} < \infty$ ;  $a = o(b)$  if  $\lim_{r \rightarrow +\infty} \frac{a(r)}{b(r)} = 0$ ;  $a = \Theta(b)$  if  $a = O(b)$  and  $b = O(a)$ .

the jobs on the server exceeds this server's resource capacity. A more general definition of  $h(\cdot)$  is given in Section 2. We assume that the resource contention does not affect the transition rates in the job model nor prompt jobs to be terminated, suitable for the application scenarios where the contention level is low and manageable. Let  $C(\sigma)$  denote the average expected cost rate per server.

Now our bin-packing problem can be formulated as follows:

$$\begin{aligned} & \underset{\sigma}{\text{minimize}} && N(\sigma) \\ & \text{subject to} && C(\sigma) \leq \epsilon, \end{aligned} \tag{1}$$

where  $\epsilon > 0$  is a budget for the cost rate of resource contention. We are interested in solving this problem in the asymptotic regime where the arrival rates  $(\Lambda_L, \Lambda_H)$  scale to infinity [33–36], motivated by the ever-increasing computing demand that drives today's computing systems to be large-scale. Specifically, we assume  $(\Lambda_L, \Lambda_H) = (\lambda_L r, \lambda_H r)$  for some fixed coefficients  $\lambda_L$  and  $\lambda_H$  and a scaling factor  $r$ , and we study the asymptotic regime where  $r$  increases.

### 1.3 Main Result

Our main result is an *asymptotically optimal* policy, named JOIN-REQUESTING-SERVER (JRS), for this new setting of SBP in service systems with time-varying job resource requirements. The asymptotic optimality is in the sense that under our proposed policy JRS, the expected number of active servers is at most  $(1 + O(r^{-0.5}))$  times the optimal objective value of the optimization problem in (1), while the cost rate incurred is at most  $(1 + O(r^{-0.5})) \cdot \epsilon$  (i.e., exceeding the budget by at most a diminishing fraction). This asymptotic optimality result translates into an *additive optimality gap* of  $O(\sqrt{r})$  in the objective value (expected number of active servers), since the optimal objective value can be shown to be  $\Theta(r)$ . This main result is formally presented in Theorem 1.

Our model can be specialized to the traditional setting of SBP in service systems where jobs' resource requirements remain fixed over time. For this specialization, we replace the constraint  $C(\sigma) \leq \epsilon$  in the problem formulation (1) with a capacity constraint, which requires the total resource requirement by jobs on a server to be within the server's resource capacity. Our proposed policy JRS can then be adapted into one that has an  $O(\sqrt{r})$  optimality gap in the objective value, which improves upon the state-of-the-art  $o(r)$  optimality gap. A discussion on the implementation complexity of JRS and how it compares with existing policies for the traditional setting of SBP in service systems is provided in Section 4.3. To be clear, this setting is not a strict special case of the formulation in (1) because  $\epsilon > 0$  is required there, but our approach and proof carry over.

From a technical approach perspective, our contribution is a novel approach that decomposes the policy design into two steps: defining a single-server sub-problem, and then converting the solution of the sub-problem into a policy in the original problem. This decomposition not only reduces the complexity of policy design but also makes the analysis tractable. We provide an overview of this approach in Section 1.4.

### 1.4 Approach Overview

To motivate our approach, we ask two questions:

*How should we design a good dispatch policy for this system?  
How can we prove that a dispatch policy is asymptotically optimal?*

Before presenting our answers to these two questions, we first comment on why they are challenging to answer. On the one hand, solving this problem directly via dynamic programming is intractable due to the unbounded state space resulting from the infinitely many servers. Even if we restrict ourselves to the servers that are active, the state space is still prohibitively large. On the other hand, we can consider designing a heuristic policy. However, unlike traditional SBP problems where we

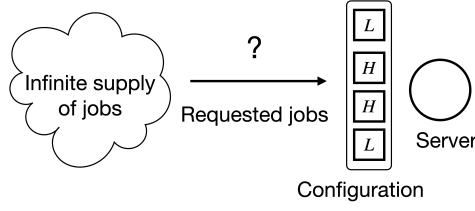


Fig. 2. A single-server system with an infinite supply of jobs. A single-server policy decides when to request jobs and how many jobs of each type to request.

can simply seek to pack the servers as compactly as possible, here the question of *how many jobs should be put on a server* is complicated. The complexity comes from the time-varying resource requirements of jobs, which affect future resource contention.

**A policy-conversion framework.** We answer the two questions at the same time with a novel *policy-conversion* framework. The framework has two steps:

- (1) Define the *single-server problem*, which is a easy-to-solve low-dimensional subproblem;
- (2) Convert the optimal policy of the single-server problem into a policy in the original problem.

This framework allows us to break down the complicated policy design problem into two components. In defining the single-server problem in Step 1, our goal is to quantify the throughput of each server under the resource contention constraint; in the policy conversion in Step 2, our goal is to dispatch jobs optimally based on each single server's characteristics. As we will show, a careful construction of the single-server problem and the conversion procedure naturally leads to a policy for the original system and a proof of its asymptotic optimality. Below we give a quick overview of how we carry out these two steps and the motivation for the design choices.

*Single-server problem.* To define the single-server problem, consider the following setting: suppose that our goal is to maximize the throughput of *one specific server* while keeping its expected cost rate of resource contention below  $\epsilon$ , then how should we send jobs to this server? Observe that even though we want to send jobs to the server as frequently as possible, the frequency is fundamentally limited by how fast the server is able to serve jobs and how many jobs can be packed on the server. This motivates us to consider the single-server system illustrated in Figure 2. The system has one server and an infinite supply of jobs of all types, so the server can start the service of any number of new jobs of any type at any time. We say the server *requests* a job from the infinite supply whenever it starts serving a new job. We assume the same job model and cost model as in the original infinite-server system. The single-server problem aims to find a job-requesting policy that maximizes the throughput (the number of each type of job that can be served) along the direction of the arrival rate vector  $(\Lambda_L, \Lambda_H) = (\lambda_L r, \lambda_H r)$ , while maintaining the steady-state expected cost rate of resource contention below  $\epsilon$ .

*How is the single-server problem related to the original problem?* Let  $\bar{N}^*$  be the number such that the total throughput of  $\bar{N}^*$  single-server systems under the optimal job-requesting policy is equal to  $(\lambda_L r, \lambda_H r)$  (assuming  $\bar{N}^*$  is an integer for simplicity). Consider the following policy in the original system: let each of the first  $\bar{N}^*$  servers in the original system adopt the optimal job-requesting policy and send requests to the dispatcher based on its current configuration. If the requested jobs were to arrive as soon as the dispatcher received the requests, the dispatcher would be able to fulfill the requests immediately. In this case, the first  $\bar{N}^*$  servers in the original system would have the same dynamics as  $\bar{N}^*$  i.i.d. single-server systems, achieving the largest possible throughput and

satisfying the constraint on resources contention. So the original system would have achieved the optimal number of active servers.

However, in the actual model, the dispatcher cannot immediately fulfill a job request because jobs arrive stochastically over time. Nevertheless, the dispatcher can still find a suitable way to match each job arrival with the requests. To see this, note that the time points when the dispatcher receives type  $i$  requests result from the superposition of  $\bar{N}^* = \Theta(r)$  independent point processes, each with the average rate  $\lambda_i r / \bar{N}^*$ . As  $r \rightarrow \infty$ , the instantaneous rates of requesting type  $i$  jobs concentrate around the arrival rate  $\lambda_i r$  for each  $i$ . As a result, most job requests can be fulfilled within a diminishing delay, so most servers in the original system can closely track the optimal single-server dynamics.

*A meta-policy, JOIN-REQUESTING-SERVER (JRS), and its asymptotic optimality.* Based on the single-server problem and the idea of tracking the optimal single-server dynamics, we propose a meta-policy, JOIN-REQUESTING-SERVER (JRS), which converts a single-server policy  $\bar{\sigma}$  to a dispatch policy in the original infinite-server system. We say that JRS takes  $\bar{\sigma}$  as a subroutine. The full definition of JRS is given in Section 4, along with discussions on various practical considerations in its implementation.

We show that the asymptotic performance of JRS is related to its subroutine in the sense described in Theorem 3, which we refer to as the *conversion theorem*. In particular, JRS with the optimal single-server policy (which we refer to as SINGLE-OPT) as the subroutine is asymptotically optimal for the original infinite-server problem as  $r \rightarrow \infty$ .

In order to track the optimal single-server dynamics, JRS uses a more sophisticated mechanism to control the long-term consequences of missing job requests or fulfilling requests with delays. The mechanism involves the auxiliary variables of *tokens* and *virtual jobs*, which regularize the process of generating requests and matching arrivals with requests. These auxiliary variables play a crucial role in the proof of Theorem 3 in Section 5, where a novel Stein's method argument is carried out. We discuss the role of tokens and virtual jobs and their necessity at the end of Section 4 and in Section 5.5.

Finally, we comment that our policy conversion framework can be applicable to other systems with similar structures. Specifically, we can try to define a suitable single-server problem, solve for its optimal policy, and convert the optimal single-server policy to the original problem. A similar conversion theorem should hold as long as the servers are weakly coupled in some sense. See Section 6 for a discussion of such systems.

*Relation to the mean-field approach.* We remark that the mean-field approach often studies the empirical distribution of configurations on all servers [17, 32–36], which can be viewed as a probability distribution of a single server's configuration. However, the mean-field approach is typically used to *analyze* this empirical distribution under a given policy for the original system. In contrast, our approach solves the single-server problem to *design* a single-server policy, and then converts it to a policy in the original system with performance guarantee.

## 1.5 Paper organization

In Section 2, we present the general problem formulation, which generalizes the simplified version in this section. In Section 3, we give a more detailed overview of our main result and approach, with a short proof of Theorem 1 (main result) based on Theorems 2–4 at the end of the section. Section 4 provides a detailed description of our meta-policy, JOIN-REQUESTING-SERVER (JRS), along with discussions on practical considerations in its implementation. In Section 5, we prove the performance guarantee of JRS (Theorem 3) under an irreducibility assumption. The proof for the



general case and other proofs are deferred to the appendices. We conclude the paper and discuss some future directions in Section 6.

## 2 PROBLEM FORMULATION

**Job Model.** As described in Section 1, we consider a job model where each job in service can be in one of multiple *phases*, each phase associated with a different resource requirement. Here the resource requirement can be a multi-dimensional vector, with each coordinate specifying the requirement of one type of resource. To model the temporal variation in the resource requirement, we assume that each job transitions between phases while in service until it is completed. The phase transition process is described by a continuous-time Markov chain on the state space  $\mathcal{I} \cup \{\perp\}$ , where  $\mathcal{I}$  is the set of phases and  $\perp$  is the absorbing state that denotes the completion of the job. We call a transition between two phases in  $\mathcal{I}$  an internal transition, and let  $\mu_{ii'}$  denote the transition rate from phase  $i$  to phase  $i'$ ; the departure of a job then corresponds to a transition from a phase  $i \in \mathcal{I}$  to  $\perp$ , whose transition rate is denoted as  $\mu_{i\perp}$ . The phase transitions of different jobs are assumed to be independent of each other.

We classify a job as a type  $i$  job if it starts from phase  $i \in \mathcal{I}$  when entering service. Jobs of each type  $i$  arrive to the system according to an independent Poisson process with rate  $\Lambda_i$ .

**Server Model.** We consider an infinite-server system with identical servers. As soon as a job arrives to the system, the job needs to be dispatched to a server to start service immediately. Note that this is always feasible because there are an infinite number of servers in the system. We assume that jobs cannot be preempted or migrated. To describe the state of a server, we define the *configuration* of a server as an  $|\mathcal{I}|$ -dimensional vector  $\mathbf{k} = (k_i)_{i \in \mathcal{I}}$ , whose  $i$ -th entry  $k_i$  is the number of jobs in phase  $i$  on the server. Each server has a limit on the total number of jobs in service at the same time. This limit is denoted as  $K_{\max}$  and referred to as the *service limit*. Then the set of feasible server configurations is  $\mathcal{K} \triangleq \{\mathbf{k} : \sum_{i \in \mathcal{I}} k_i \leq K_{\max}\}$ .

**System Dynamics.** The system state can be represented by the concatenation of the configurations of all servers. Specifically, we index the servers by positive integer numbers and denote the configuration of server  $\ell$  at time  $t$  as  $\mathbf{K}^\ell(t)$ . Then the state of the entire system can be represented by the infinite vector  $(\mathbf{K}^\ell(t))_{\ell \in \mathbb{Z}_+}$ .

Suppose that the system is in state  $(\mathbf{k}^\ell)_{\ell \in \mathbb{Z}_+}$ . Let  $\mathbf{e}_i$  be an  $|\mathcal{I}|$ -dimensional vector whose  $i$ -th entry is 1 and all other entries are 0. Then the following state transitions can happen:

- $\mathbf{k}^\ell \rightarrow \mathbf{k}^\ell + \mathbf{e}_i, \mathbf{k}^{\ell'} \rightarrow \mathbf{k}^{\ell'} \forall \ell' \neq \ell$ : a type  $i$  job arrives and is dispatched to server  $\ell$ ;
- $\mathbf{k}^\ell \rightarrow \mathbf{k}^\ell + \mathbf{e}_{i'} - \mathbf{e}_i, \mathbf{k}^{\ell'} \rightarrow \mathbf{k}^{\ell'} \forall \ell' \neq \ell$ : a job on server  $\ell$  transits from phase  $i$  to phase  $i'$ ;
- $\mathbf{k}^\ell \rightarrow \mathbf{k}^\ell - \mathbf{e}_i, \mathbf{k}^{\ell'} \rightarrow \mathbf{k}^{\ell'} \forall \ell' \neq \ell$ : a job on server  $\ell$  departs the system from phase  $i$ .

The specifics of the system dynamics depend on the employed *dispatch policy* that decides which server to dispatch to when a job arrives.

**Active Servers.** We are interested in the number of *active servers*, i.e., servers currently serving a positive number of jobs. Note that given the arrival rates of jobs, the smaller the number of active servers, the better the system is utilized. Let  $X_{\mathbf{k}}(t)$  be the number of servers in configuration  $\mathbf{k}$  at time  $t$ , i.e.,  $X_{\mathbf{k}}(t) = \sum_{\ell=1}^{\infty} \mathbb{1}_{\{\mathbf{K}^\ell(t)=\mathbf{k}\}}$ . Then the number of active servers can be written as  $\sum_{\mathbf{k} \neq \mathbf{0}} X_{\mathbf{k}}(t)$ , where  $\mathbf{0} \in \mathbb{R}^{|\mathcal{I}|}$  is the zero vector.

**Cost of Resource Contention.** Recall that the cost rate function  $h(\cdot)$  maps a server's configuration to a rate of cost. We assume that  $h(\cdot)$  is any function that is  $\Gamma$ -Lipschitz continuous with respect to the  $L^1$  distance for some constant  $\Gamma > 0$  and satisfies  $h(\mathbf{0}) = 0$ .

**Performance Goal.** Our high-level goal is to design dispatch policies that minimize the number of active users while keeping the cost rate of resource contention within a certain budget. Specifically, we consider policies that are allowed to be randomized and non-Markovian (i.e., the policies can make history-dependent decisions). We further focus on policies that induce a unique stationary distribution on the configuration process  $\{(K^\ell(t))_{t \in \mathbb{Z}_+}\}$ , assuming that the configuration process is embedded in a Markov chain that has a unique stationary distribution. We are interested in such policies because the resulting time averages of quantities related to the configurations are equal to the corresponding expectations under the unique stationary distribution regardless of the initial state. Let  $\sigma$  be a policy of interest,  $(K^\ell)_{\ell \in \mathbb{Z}_+}$  be a random element that follows the stationary distribution of the system state induced by  $\sigma$ , and  $X_k$  be the corresponding number of servers in configuration  $k$  in steady state under  $\sigma$ . Then the expected number of active servers is given by

$$N(\sigma) \triangleq \sum_{k \neq 0} \mathbb{E}[X_k].$$

We define the expected cost rate per expected active server as

$$C(\sigma) \triangleq \frac{\sum_{k \neq 0} h(k) \mathbb{E}[X_k]}{\sum_{k \neq 0} \mathbb{E}[X_k]}.$$

Note that if  $C(\sigma) \leq \epsilon$ , we have  $\sum_{k \neq 0} h(k) \mathbb{E}[X_k] \leq \epsilon \sum_{k \neq 0} \mathbb{E}[X_k]$ . Now our goal can be formulated as the following optimization problem, referred to as problem  $\mathcal{P}((\Lambda_i)_{i \in \mathcal{I}}, \epsilon)$ :

$$\begin{aligned} & \underset{\sigma}{\text{minimize}} && N(\sigma) \\ & \text{subject to} && C(\sigma) \leq \epsilon, \end{aligned} \tag{2}$$

where  $\epsilon$  is a budget for the cost rate of resource contention.

**Asymptotic Optimality.** We focus on the asymptotic regime where for all  $i \in \mathcal{I}$ , the arrival rate is given by  $\Lambda_i = \lambda_i r$  for some constant coefficient  $\lambda_i$  and a positive *scaling factor*  $r \rightarrow +\infty$ . To define asymptotic optimality, we first define the following notion of approximation to the optimization problem  $\mathcal{P}((\Lambda_i)_{i \in \mathcal{I}}, \epsilon)$  in (2): a policy  $\sigma$  is said to be  $(\alpha, \beta)$ -optimal if  $N(\sigma) \leq \alpha \cdot N^*((\Lambda_i)_{i \in \mathcal{I}}, \epsilon)$  and  $C(\sigma) \leq \beta \cdot \epsilon$ , where  $N^*((\Lambda_i)_{i \in \mathcal{I}}, \epsilon)$  is the optimal objective value in (2). Now consider a family of policies  $\sigma^{(r)}$  indexed by the scaling factor  $r$ . We say that the policy  $\sigma^{(r)}$  is *asymptotically optimal* if it is  $(\alpha^{(r)}, \beta^{(r)})$ -optimal to the optimization problem  $\mathcal{P}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$  with  $\alpha^{(r)}, \beta^{(r)} \rightarrow 1$  as  $r \rightarrow \infty$ . We will suppress the superscript  $(r)$  for simplicity when it is clear from the context.

We note that under any policy  $\sigma$ ,  $N(\sigma) = \Theta(r)$ . This can be proven using the renowned Little's Law [19] in the following way. The total job arrival rate is  $\Theta(r)$  and the expected time that a job spends in the system is  $O(1)$ . So by Little's Law, the expected number of jobs in the system in steady state is  $\Theta(r)$ . Since each server can accommodate a constant number of jobs, the expected number of active servers is  $\Theta(r)$ . Given this,  $(\alpha^{(r)}, \beta^{(r)})$ -optimality implies an optimality gap of  $\alpha^{(r)} \cdot r$  in the objective value.

### 3 MAIN RESULT AND OUR APPROACH

#### 3.1 Main Result

Our main result, Theorem 1, is the asymptotic optimality of our proposed policy JOIN-REQUESTING-SERVER (JRS), with a subroutine we call SINGLE-OPT, as briefly discussed in Section 1. This asymptotic optimality result implies an  $O(\sqrt{r})$  optimality gap in the expected number of active servers. We defer the detailed descriptions of JRS and SINGLE-OPT to Section 4 and Appendix C. Theorem 1 follows immediately from Theorems 2–4 to be introduced in Section 3.2; a short proof is included at the end of Section 3.2 for clarity.



**Theorem 1** (Asymptotic Optimality). *Consider a stochastic bin-packing problem in service systems with time-varying job resource requirements. Let the arrival rates be  $(\lambda_i r)_{i \in \mathcal{I}}$  and the cost rate budget be  $\epsilon > 0$ . Then the policy JOIN-REQUESTING-SERVER (JRS) with the subroutine SINGLE-OPT is  $(1 + O(r^{-0.5}), 1 + O(r^{-0.5}))$ -optimal. That is, the expected number of active servers under JRS with SINGLE-OPT is at most  $(1 + O(r^{-0.5}))$  times the optimal value of the problem  $\mathcal{P}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$ , while the cost rate incurred is at most  $(1 + O(r^{-0.5})) \cdot \epsilon$ .*

**Specialization to Non-Time-Varying Resource Requirements.** As mentioned in Section 1, we can specialize this result to the setting where the resource requirement of a job does not vary over time. To do that, we remove the cost constraint in  $\mathcal{P}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$ , and redefine the set of feasible server configurations,  $\mathcal{K}$ , to also incorporate hard capacity constraints for each type of resources. The rest of the analysis is almost identical to that of the analysis for time-varying resource requirements; we omit the details due to the space limit. This specialization results in a policy that is  $(1 + O(r^{-0.5}))$ -optimal in the expected number of active servers.

### 3.2 Our Approach

In a nutshell, our approach is to reduce the original optimization problem in an infinite-server system to an optimization problem in a single-server system, which is defined below.

**A Single-Server System.** Consider a single-server system serving jobs with time-varying resource requirements. The system has an infinite supply of jobs of all types. As a result, the server can request any number of new jobs of any type at any time. Once a job is requested, it immediately enters service.

We represent the server configuration at time  $t$  using a vector  $\bar{\mathbf{K}}(t) = (\bar{K}_i(t))_{i \in \mathcal{I}}$ , whose  $i$ -th entry denotes the number of jobs in phase  $i$ . We assume that the single-server system has the same service limit  $K_{\max}$  and cost rate function  $h(\cdot)$  as a server in the original infinite-server system. Therefore, the server configuration  $\bar{\mathbf{K}}(t)$  is also in the set  $\mathcal{K} = \{\mathbf{k} : \sum_{i \in \mathcal{I}} k_i \leq K_{\max}\}$ , and the cost rate at time  $t$  is  $h(\bar{\mathbf{K}}(t))$ .

A single-server policy  $\bar{\sigma}$  determines when and how many jobs of each type to request. We allow the single-server policy to be randomized and assume it is Markovian, i.e., it makes decisions only based on the current configuration. Note that allowing non-Markovian policies will not change the optimal value of the single-server problem that we will consider (see Appendix C). Let  $\pi \triangleq (\pi(\mathbf{k}))_{\mathbf{k} \in \mathcal{K}}$  be a stationary distribution of the server configuration under the policy  $\bar{\sigma}$ , and let  $\bar{\mathbf{K}}(\infty)$  be a random variable with the distribution  $\pi$ . When we consider a policy  $\bar{\sigma}$  and its stationary distribution  $\pi$ , we assume that the system is initialized from  $\pi$ . The policy  $\bar{\sigma}$  together with  $\pi$  defines the request rate of type  $i$  jobs  $\bar{\lambda}_i$ , which is the expected number of type  $i$  jobs requested per unit time in steady state. Note that  $\bar{\lambda}_i$  is the throughput of type  $i$  jobs since the system has a finite state space.

We consider the following single-server problem, denoted as  $\bar{\mathcal{P}}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$ :

$$\begin{aligned} & \underset{\bar{N}, \bar{\sigma}, \pi}{\text{minimize}} && \bar{N} \\ & \text{subject to} && \mathbb{E}[h(\bar{\mathbf{K}}(\infty)) | \bar{\mathbf{K}}(\infty) \neq \mathbf{0}] \leq \epsilon, \\ & && \bar{N} \bar{\lambda}_i = \lambda_i r, \quad \forall i \in \mathcal{I}. \end{aligned} \tag{3}$$

The single-server problem can be interpreted as follows. We can think of  $\bar{N}$  as the number of copies of the single-server system under  $\bar{\sigma}$  needed to support the arrival rates  $(\lambda_i r)_{i \in \mathcal{I}}$  in the infinite-server system. To minimize  $\bar{N}$ , it is equivalent to maximizing the throughput  $(\bar{\lambda}_i)_{i \in \mathcal{I}}$  in each single-server system, while maintaining their proportions as  $(\lambda_i r)_{i \in \mathcal{I}}$ .

We remark that for the problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$ , we only need to consider policies that do not depend on the scaling factor  $r$ . To see this, we can replace the decision variable  $\overline{N}$  with  $\bar{n} \triangleq \overline{N}/r$  and the optimization problem can be equivalently formulated as follows, which does not involve  $r$ :

$$\begin{aligned} & \underset{\bar{n}, \bar{\sigma}, \pi}{\text{minimize}} && \bar{n} \\ & \text{subject to} && \mathbb{E} [h(\overline{K}(\infty)) | \overline{K}(\infty) \neq \mathbf{0}] \leq \epsilon, \\ & && \bar{n} \bar{\lambda}_i = \lambda_i, \quad \forall i \in I. \end{aligned} \quad (4)$$

**Lower Bound.** The single-server problem gives a lower bound to the original problem in (2) as stated in the following theorem. The proof is given in Appendix A.

**Theorem 2 (Lower Bound).** *Consider a stochastic bin-packing problem in service systems with time-varying job resource requirements. Let the arrival rates be  $(\lambda_i r)_{i \in I}$  and the cost rate budget be  $\epsilon > 0$ . Let  $N^*$  be the optimal value of the original infinite-server problem in (2), and let  $\overline{N}^*$  be the optimal value of the single-server problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$ , then  $N^* \geq \overline{N}^*$ .*

**Converting From the Single-Server System to the Infinite-Server System.** Having established a lower bound on the infinite-server problem  $\mathcal{P}((\lambda_i r)_{i \in I}, \epsilon)$  in terms of the optimal value of the single-server problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$ , next we focus on finding an asymptotically optimal policy. We will characterize the performance guarantee of a class of policies and then show that the best policy within the class is asymptotically optimal. Specifically, we consider a meta-policy called JOIN-REQUESTING-SERVER (JRS), which converts a Markovian single-server policy  $\bar{\sigma}$  into an infinite-server policy. We call the policy resulting from the conversion a *JRS policy with a subroutine  $\bar{\sigma}$* . Through analyzing the meta-policy JRS, we show that the performance of each JRS policy can be characterized by the performance of its subroutine, as stated in Theorem 3 below. The proof of Theorem 3 under an irreducibility assumption is given in Section 5, and the proof for the full version is given in Appendix B.3.

**Theorem 3 (Conversion Theorem).** *Consider a stochastic bin-packing problem in service systems with time-varying job resource requirements. Let the arrival rates be  $(\lambda_i r)_{i \in I}$  and the cost rate budget be  $\epsilon > 0$ . Let  $(\overline{N}, \bar{\sigma}, \pi(\mathbf{k}))$  be a solution feasible to the single-server problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$ . In addition, we assume that the policy  $\bar{\sigma}$  is Markovian. Let the infinite-server policy  $\sigma$  be JRS with a subroutine  $\bar{\sigma}$ . Then under  $\sigma$ , we have*

$$\left| \sum_{k \neq \mathbf{0}} \mathbb{E} [X_k] - \lceil \overline{N} \rceil \cdot \mathbb{P}(\overline{K} \neq \mathbf{0}) \right| = O(\sqrt{r}), \quad (5)$$

$$\left| \sum_{k \neq \mathbf{0}} h(\mathbf{k}) \mathbb{E} [X_k] - \lceil \overline{N} \rceil \cdot \mathbb{E} [h(\overline{K})] \right| = O(\sqrt{r}). \quad (6)$$

As a result,

$$N(\sigma) \leq (1 + O(r^{-0.5})) \cdot \overline{N}, \quad (7)$$

$$C(\sigma) \leq (1 + O(r^{-0.5})) \cdot \epsilon. \quad (8)$$

**Optimal Single-Server Policy.** Theorem 3 together with the lower bound in Theorem 2 reduces the infinite-server problem  $\mathcal{P}((\lambda_i r)_{i \in I}, \epsilon)$  in (2) to the single-server problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$  in (3). We can obtain the optimal single-server policy, SINGLE-OPT, by solving a linear program, as stated in the theorem below.

**Theorem 4** (Optimality of Single-OPT, Informal). *There exists a linear program  $\overline{\mathcal{LP}}((\lambda_i)_{i \in I}, \epsilon)$  that is equivalent to the single-server problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$ . In particular, we can construct an optimal Markovian policy for  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$  from the optimal solution of  $\overline{\mathcal{LP}}((\lambda_i)_{i \in I}, \epsilon)$ .*

Proof of Theorem 4 and details on the construction of the optimal policy are given in Appendix C.

**PROOF OF THE THEOREM 1 BASED ON THEOREMS 2–4.** Because SINGLE-OPT (along with an optimal stationary distribution) optimally solves  $\overline{\mathcal{P}}((\lambda_i r)_{i \in I}, \epsilon)$ , it achieves the optimal value  $\overline{N}^*$ . Let  $\sigma$  be JRS with a subroutine SINGLE-OPT, then according to Theorem 3, we have  $N(\sigma) \leq (1 + O(r^{-0.5})) \cdot \overline{N}^*$  and  $C(\sigma) \leq (1 + O(r^{-0.5})) \cdot \epsilon$ . By Theorem 2, we also have  $N^* \geq \overline{N}^*$ . So we conclude that JRS with a subroutine SINGLE-OPT is  $(1 + O(r^{-0.5}), 1 + O(r^{-0.5}))$ -optimal.  $\square$

## 4 PROPOSED META-POLICY: JOIN-REQUESTING-SERVER (JRS)

In this section, we describe our meta-policy, JOIN-REQUESTING-SERVER (JRS), in full detail. For ease of presentation, we focus on the case where the subroutine policy  $\bar{\sigma}$  for JRS is  $\mathbf{k}^0$ -irreducible, i.e., under  $\bar{\sigma}$ , there exists a configuration  $\mathbf{k}^0$  such that the single-server system can return to  $\mathbf{k}^0$  from any other configurations (which is equivalent to assuming that the configuration of the single-server system under policy  $\bar{\sigma}$  forms a *unichain*). The algorithm for the general case is given in Appendix B.3.

### 4.1 How the Single-Server Policy Requests Jobs

Before going into the definition of JRS, we first take a closer look at how the Markovian *single-server* policy  $\bar{\sigma}$  requests jobs, to avoid potential ambiguity caused by the fact that a single-server policy can request jobs at *any time*. Let  $a_i$  denote the number of type  $i$  jobs requested, and let  $\mathbf{a} \triangleq (a_i)_{i \in I}$ . We say  $\mathbf{a}$  is *feasible* if the total number of jobs on the server does not exceed  $K_{\max}$  after adding the jobs. The policy  $\bar{\sigma}$  performs one of the following two types of requests based on the current configuration.

- **Reactive requests.** A *reactive request* is triggered by *either an internal transition or a departure*. The changes in the configuration when a reactive request is made can be represented by the diagram

$$\mathbf{k} \rightarrow \mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a},$$

where  $\mathbf{k} \rightarrow \mathbf{k}'$  is due to the internal transition or departure, and  $\mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a}$  happens since the policy immediately requests  $\mathbf{a}$  jobs. The policy  $\bar{\sigma}$  specifies a probability distribution over all feasible  $\mathbf{a}$  when it decides to perform reactive requests for the configuration  $\mathbf{k}'$ .

- **Proactive requests.** A *proactive request* happens on its own, and it happens at a *finite rate* depending on the current configuration of the server. The change of the configuration when a proactive request happens can be represented by the diagram

$$\mathbf{k} \rightarrow \mathbf{k} + \mathbf{a}.$$

More specifically, suppose the policy  $\bar{\sigma}$  decides to perform proactive requests for a configuration  $\mathbf{k}$ . Then for each feasible  $\mathbf{a}$ , the policy  $\bar{\sigma}$  specifies a rate and runs a timer with an exponentially distributed duration with the specified rate. When the timer ticks,  $\mathbf{a}$  jobs are requested. When the configuration changes, all the timers are canceled and restarted with new rates based on the new configuration.

### 4.2 Description of JOIN-REQUESTING-SERVER (JRS)

The inputs of JRS include: (i) the single-server policy  $\bar{\sigma}$ , (ii) the objective value of  $\bar{\sigma}$  in the single-server problem (3), denoted as  $\overline{N}$ , and (iii) the transition rates in the job model.

We first divide the infinite server pool into two sets based on the server index  $\ell$ . Let  $L = \lceil \bar{N} \rceil$ . We call servers with index  $\ell \leq L$  *normal servers*; we call servers with index  $\ell > L$  *backup servers*. The normal servers are responsible for serving most of the jobs, while the backup servers are activated only to handle overflow jobs (jobs that are not dispatched to normal servers).

The JRS is specified in two steps.

- **Step 1 (Job Requesting on a Normal Server):** We let each normal server request jobs using its subroutine, the single-server policy  $\bar{\sigma}$ . The input to the policy  $\bar{\sigma}$  is what we refer to as the *observed configuration* of the server, which will be further explained below. When  $\bar{\sigma}$  requests  $\mathbf{a} = (a_i)_{i \in \mathcal{I}}$  jobs,  $a_i$  type  $i$  *tokens* are generated for each  $i \in \mathcal{I}$  to store the job requests. The server *pauses* the job requesting process if it already has any type of tokens, and resumes when all the tokens that it generated are removed.
- **Step 2 (Arrival Dispatching):**
  - **Real jobs.** When a type  $i$  job arrives, the dispatcher chooses a type  $i$  token uniformly at random, removes the token, and assigns the job to the corresponding server. When there are no type  $i$  tokens, the dispatcher sends the job to an idle backup server.
  - **Virtual jobs.** When the total number of type  $i$  tokens throughout the system exceeds the limit  $\eta_{\max} = \lceil \sqrt{L} \rceil$  (called the *token limit*), a type  $i$  *virtual arrival* is triggered, which causes the dispatcher to choose a type  $i$  token uniformly at random, remove the token, and assign a *virtual job* to the corresponding server. A virtual job has the same transition dynamics as a real job but does not consume physical resources.

The *observed configuration* of a normal server in Step 1 is the configuration resulting from real jobs and virtual jobs combined. That is, it is a vector whose  $i$ -th entry represents the total number of real and virtual jobs in phase  $i$  on this server. The observed configuration changes when there is a new real or virtual job arrival assigned to the server, or when a real or virtual job on the server has a phase transition or departs. We update the input to the policy  $\bar{\sigma}$  when the observed configuration changes. Whenever the observed configuration changes, the policy  $\bar{\sigma}$  cancels the exponential timers in progress; but a reactive request from the policy  $\bar{\sigma}$  can only be triggered when a real or virtual job on the server has a phase transition or departs.

**Intuition behind JRS.** To provide a better understanding of the main design ideas of JRS, here we give an intuitive description of how it works. Broadly, servers generate job requests and store unfulfilled requests as tokens; the dispatcher assigns jobs to servers according to the tokens to fulfill job requests. This is the mechanism for matching job arrivals with requests, which is referenced at the end of Section 1.4. However, rather than matching all tokens with job arrivals, JRS opts to convert some of the tokens into virtual jobs to keep the total number of tokens within an upper limit  $\eta_{\max}$ . By capping the number of tokens, JRS ensures that the job requests generated by each server get fulfilled quickly (either by a real job or a virtual job), and thus the observed configurations of servers maintain proximity to i.i.d. copies of the single-server systems.

The choice of the token limit  $\eta_{\max} = \Theta(\sqrt{r})$  balances two key considerations. On the one hand, a smaller  $\eta_{\max}$  brings the observed configurations closer to i.i.d. copies of single-server systems. On the other hand, if  $\eta_{\max}$  is overly small, the rate of generating virtual jobs becomes high and the probability for a job arrival to see no tokens is also high. As a result, the observed configurations, which include both real and virtual jobs, deviate from the real-job configurations. A more in-depth discussion on the role of tokens and virtual jobs and whether they are fundamental is in Section 5.5.

### 4.3 Practical considerations in implementing JOIN-REQUESTING-SERVER (JRS)

**Computational complexity of JRS.** The computational complexity of JRS consists of two components: the *offline* component that computes a single-server policy  $\bar{\sigma}$  and its objective value  $\bar{N}$ , and the *online* component that carries out the two steps of JRS.

The offline component reduces to solving the linear program given in (61) in Appendix C.1, whose number of optimization variables is linear in the number of feasible configurations times the number of phases, i.e.,  $|\mathcal{K}| \times |I|$ , on a single server. Admittedly,  $|\mathcal{K}|$  can be large when a single server has a large quantity of resources and there are many job phases. However, we opt for the view that a single server is not excessively large and the system's scale is primarily captured by the scaling factor  $r$ . Therefore, it is advantageous that the computational complexity of this offline component is independent of  $r$ .

In the online component, the bulk of the computation is in job requesting and virtual job simulation, which can be executed distributedly on the normal servers. Specifically, each normal server monitors its observed configuration and generates tokens according to the single-server policy  $\bar{\sigma}$ ; additionally, when a virtual job is assigned to the server, the server simulates the dynamics of the virtual job, i.e., generates random variables corresponding to phase transitions and job departure. Backup servers do not need to perform any computation beyond serving jobs.

The scheduler, which stores all the tokens, has two responsibilities in the online component: (i) the scheduler matches each newly arrived job to a token of the same type, chosen uniformly at random, or sends the job to a backup server when there are no tokens of the same type; (ii) the scheduler monitors the number of tokens of each type and assigns virtual jobs when the number of tokens exceeds the limit  $\eta_{\max}$ .

It is informative to compare the computational complexity of JRS with existing algorithms designed for the traditional setting of stochastic bin-packing in service systems, where the resource requirements are non-time-varying [17, 32–36]. At a high level, these existing algorithms function as follows: upon the arrival of a job, the scheduler checks the current configurations of all servers and assigns the job to a server whose configuration optimizes certain predefined criteria. Among these, the GRAND algorithm [33, 35, 36] stands out for its simplicity and asymptotic optimality. Under GRAND, the scheduler only needs to identify configurations that can accommodate the incoming job and then randomly assigns the job to one of these feasible servers, along with some idle servers. Compared to JRS, GRAND does not have an offline planning component, and individual servers do not perform computation beyond serving jobs. The scheduler's role in GRAND is slightly more complex than in JRS. Consequently, when considering using JRS in settings where job resource requirements are non-time-varying, practitioners should weigh whether the additional computational complexity is warranted.

**Model parameter estimation.** A limitation of JRS is its dependency on known model parameters, including job arrival rates and phase transition rates. Such dependency is not present in existing algorithms designed for the setting with non-time-varying resource requirements. In real-world applications, the model parameters can be estimated from workload traces such as [37, 38]. Estimation errors can impact the system's performance, an issue that merits further in-depth investigation in future work. Here, we provide a preliminary result on the performance degradation due to parameter estimation errors. Roughly speaking, suppose that the estimation error in the job arrival rate coefficients  $\lambda_i$ 's and the phase transition rates  $\mu_{i'}$ 's and  $\mu_{i\perp}$ 's are bounded by  $\delta \geq 0$  (along with an insensitivity assumption on the single-server problem). Then if we use JRS where the single-server policy is obtained by solving for the optimal single-server policy under the estimated parameters, the resulting JRS is  $(1 + B\delta + O(r^{-0.5}), 1 + B\delta + O(r^{-0.5}))$ -optimal for any  $\delta \leq \delta_{\max}$ , where  $B$  and

$\delta_{\max}$  are positive constants independent of  $r$ . The exact statement is given in Proposition 1 in Appendix D, along with a proof.

**Connection to practical algorithms.** Recent progress has been made in addressing the issue of low utilization due to time-varying job resource requirements, notably within Google's datacenters, as discussed in [2]. The approach in [2] makes predictions on the future resource requirements of jobs, which lead to a further prediction on the future peak resource requirement on a server if a newly arrived job were to be sent to that server (assuming no future job arrivals). This prediction categorizes each server as either feasible or infeasible for the new job, and this binary outcome is subsequently used by a separate scheduler for job assignment.

Our proposed JRS policy can be viewed as giving more detailed predictions on whether it is suitable for a server to take on new jobs, represented by the tokens. The predictions are optimized by taking into account future job arrivals and the stochastic dynamics of jobs.

## 5 PROOF OF THEOREM 3 (CONVERSION THEOREM) ASSUMING IRREDUCIBILITY

In this section, we prove Theorem 3 to establish the performance guarantee of JRS. For ease of presentation, we focus on the case where the subroutine policy  $\bar{\sigma}$  is  $k^0$ -irreducible. The proof for the general case is in Appendix B.3.

This section is organized as follows. We first provide some preliminaries in Section 5.1. Then we outline the steps and lemmas needed for the proof in Section 5.2. In Section 5.3, we prove Theorem 3 based on the lemmas. In Section 5.4, we prove one of the lemmas, Lemma 2, where we devise a novel approach to employ Stein's method. Finally, in Section 5.5, we discuss the role of tokens and virtual jobs and their necessity from a proof perspective. The proofs of the rest of the lemmas presented in this section are given in Appendix B.

### 5.1 Preliminaries

Consider an infinite-server system under the JRS policy. For each normal server  $\ell$ , we describe its status at time  $t$  using the following variables: *configuration of real jobs*  $K^\ell(t)$  (referred to simply as configuration in previous sections), *tokens*  $\eta^\ell(t)$ , *configuration of virtual jobs*  $\zeta^\ell(t)$ , and *observed configuration*  $\widehat{K}^\ell(t) \triangleq K^\ell(t) + \zeta^\ell(t)$ . We use the superscript " $1:L$ " to refer to a certain descriptor of all normal servers, for example,  $\widehat{K}^{1:L}(t) \triangleq (\widehat{K}^\ell(t))_{\ell=1,2,\dots,L}$ . The system under JRS is a Markov chain with a unique Markovian representation  $((K^\ell(t))_{\ell=1,2,\dots,L}, \zeta^{1:L}(t), \eta^{1:L}(t))$ . The following lemma shows that the system has a unique stationary distribution (the proof is provided in Appendix B.1).

**Lemma 1 (Unique Stationary Distribution).** *Consider an infinite-server system under the JRS policy with  $\bar{\sigma}$  as its subroutine, where  $\bar{\sigma}$  is a single-server policy that is Markovian and  $k^0$ -irreducible. Then the state of the system  $((K^\ell(t))_{\ell=1,2,\dots,L}, \zeta^{1:L}(t), \eta^{1:L}(t))$  has a unique stationary distribution.*

Let  $\bar{K}^{1:L}(t) \triangleq (\bar{K}^\ell(t))_{\ell=1,2,\dots,L}$  be the configuration vector of  $L$  i.i.d. copies of the single-server system under  $\bar{\sigma}$ . As discussed in Section 4.2, we will show that  $K^{1:L}(\infty)$  can be approximated by  $\bar{K}^{1:L}(\infty)$ . In the remainder of this section, we omit the steady-state symbol  $(\infty)$  for simplicity.

To rigorously discuss the approximation of the steady-state random variables, we define some metrics. Recall that  $\mathcal{K} \triangleq \{k: \sum_{i \in \mathcal{I}} k_i \leq K_{\max}\}$  is the set of feasible single-server configurations. Let  $\mathcal{K}^L \triangleq \{k^{1:L}: k^\ell \in \mathcal{K}, \forall \ell\}$  be the set of feasible configurations for all normal servers. We use  $\|\cdot\|$  to denote the  $L^1$  norm in both space  $\mathcal{K}$  and space  $\mathcal{K}^L$ :

$$\begin{aligned} \|k - k'\| &= \sum_{i \in \mathcal{I}} |k_i - k'_i|, \quad \text{for } k, k' \in \mathcal{K}, \\ \|k^{1:L} - k'^{1:L}\| &= \sum_{\ell=1}^L \|k^\ell - k'^\ell\|, \quad \text{for } k^{1:L}, k'^{1:L} \in \mathcal{K}^L. \end{aligned}$$



For any two random variables  $U^a, U^b \in \mathcal{K}^L$ , their closeness will be measured in terms of Wasserstein distance as follows:

$$d(U^a, U^b) \triangleq \sup_{f \in \text{Lip}(1)} \left\{ \mathbb{E}[f(U^a)] - \mathbb{E}[f(U^b)] \right\},$$

where the supremum is taken over the all Lipschitz-1 functions from  $\mathcal{K}^L$  to  $\mathbb{R}$ .

## 5.2 Steps and Lemmas Needed for the Proof of Theorem 3 Assuming Irreducibility

Our goal is to show that the steady-state distribution of the normal servers' real-job configurations  $K^{1:L}$  is close to the steady-state distribution of i.i.d. copies of the single-server systems  $\bar{K}^{1:L}$  in Wasserstein distance, and that the backup servers are almost empty as the arrival rate gets large. More formally, we aim to show that  $d(K^{1:L}, \bar{K}^{1:L}) = O(\sqrt{r})$  and  $\sum_{\ell=L+1}^{\infty} \sum_{i \in I} K_i^\ell = O(\sqrt{r})$  as  $r \rightarrow \infty$ . These two bounds provide the performance guarantee claimed in Theorem 3.

Instead of directly looking into the distribution of real-job configuration  $K^{1:L}$ , we show that the distribution of each of the three sums,  $K^{1:L} + \zeta^{1:L} + \eta^{1:L}$ ,  $K^{1:L} + \zeta^{1:L}$ , and  $K^{1:L}$ , can be approximated by the distribution of  $\bar{K}^{1:L}$  in Wasserstein distance. The approximation result for each sum helps us derive the approximation result for the sum with one fewer term. The result that the backup servers are almost empty also follows from these approximations. This sequence of approximations is illustrated in the figure below, where recall that  $\hat{K}^\ell(t) \triangleq K^\ell(t) + \zeta^\ell(t)$  is the observed configuration.

$$\hat{K}^{1:L} \approx \boxed{K^{1:L} + \zeta^{1:L} + \eta^{1:L}} \approx \bar{K}^{1:L}$$

A crucial observation that leads to this stepwise proof is that the process  $(\hat{K}^{1:L}(t), \eta^{1:L}(t))$  forms a Markov chain on its own. This is because real jobs and virtual jobs have the same transition dynamics and are indistinguishable by the subroutine when requesting jobs. Moreover, by the construction of JRS, the Markov chain  $(\hat{K}^{1:L}(t), \eta^{1:L}(t))$  governs the dynamics of the virtual-job configurations  $\zeta^{1:L}(t)$  and the configurations on backup servers.

Our proof consists of two steps. In **Step 1**, we focus on the process  $(\hat{K}^{1:L}(t), \eta^{1:L}(t))$ . We show that  $d(\hat{K}^{1:L} + \eta^{1:L}, \bar{K}^{1:L}) = O(\sqrt{r})$ , which immediately implies  $d(\hat{K}^{1:L}, \bar{K}^{1:L}) = O(\sqrt{r})$  because we have limited the total number of tokens to  $O(\sqrt{r})$ . In **Step 2**, we use the approximation result for  $\hat{K}^{1:L}$  in **Step 1** to show that the total number of virtual jobs,  $\sum_{i \in I} \sum_{\ell=1}^L \zeta_i^\ell$ , and the total number of jobs on backup servers are both  $O(\sqrt{r})$ . Recall that  $K^{1:L} = \hat{K}^{1:L} - \zeta^{1:L}$ , so we get  $d(K^{1:L}, \bar{K}^{1:L}) = O(\sqrt{r})$ .

Next, we state the specific lemmas.

**Step 1.** Lemma 2 below bounds the Wasserstein distance between  $\hat{K}^{1:L}$  and  $\bar{K}^{1:L}$ .

**Lemma 2.** Under the conditions of Theorem 3 and  $\bar{\sigma}$  being  $k^0$ -irreducible, we have

$$d(\hat{K}^{1:L}, \bar{K}^{1:L}) = O(\sqrt{r}).$$

The key challenge for proving Lemma 2 is that the job dispatching decisions are based on the configurations of *all normal servers*, which creates dependencies among the transitions of different servers. The key idea that helps us overcome this challenge is to consider the sum  $\hat{K}^{1:L} + \eta^{1:L}$ , which remains unchanged under job arrivals regardless of dispatching decisions. Observe that  $\hat{K}^{1:L} + \eta^{1:L}$  has *decoupled* dynamics across servers because it is only changed by internal phase transitions, departures, and requests of new jobs, which happen independently on each server. This helps us prove  $d(\hat{K}^{1:L} + \eta^{1:L}, \bar{K}^{1:L}) = O(\sqrt{r})$ , which implies Lemma 2, as argued earlier in the section.

Formally, the proof of Lemma 2 makes use of Stein's method (see, e.g., [3–5]) to compare  $\hat{K}^{1:L} + \eta^{1:L}$  with  $\bar{K}^{1:L}$ . Stein's method usually consists of three steps: generator comparison, Stein factor bound,

and moment bound. In our case, due to the finiteness of the state space  $\mathcal{K}$ , we only need to do the generator comparison and the Stein factor bound. In the generator comparison step, we show that the instantaneous transition rates of  $\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}$  match with those of  $\overline{\mathbf{K}}^{1:L}$ ; in the Stein factor bound step, we show that small difference in the transition rates of  $\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}$  and  $\overline{\mathbf{K}}^{1:L}$  does not cause much increase in the overall distance of the distributions. The detailed proof is in Section 5.4.

**Step 2.** We establish Lemma 3 and Lemma 4 below, which bound the steady-state expected number of virtual jobs and the jobs on backup servers.

**Lemma 3.** *Under the conditions of Theorem 3 and  $\bar{\sigma}$  being  $\mathbf{k}^0$ -irreducible, for each  $i \in \mathcal{I}$ , the steady-state expected number of virtual jobs of type  $i$  is of the order  $O(\sqrt{r})$ , i.e.,*

$$\mathbb{E} \left[ \sum_{\ell=1}^L \zeta_i^\ell \right] = O \left( \sqrt{r} \right).$$

**Lemma 4.** *Under the conditions of Theorem 3 and  $\bar{\sigma}$  being  $\mathbf{k}^0$ -irreducible, for each  $i \in \mathcal{I}$ , the steady-state expected number of type  $i$  jobs on backup servers is of the order  $O(\sqrt{r})$ , i.e.,*

$$\mathbb{E} \left[ \sum_{\ell=L+1}^\infty K_i^\ell \right] = O \left( \sqrt{r} \right).$$

The key idea for proving Lemma 3 and Lemma 4 is that by the characterization of  $\widehat{\mathbf{K}}^{1:L}$  in Lemma 2 and the fact that the job requests are made based on  $\widehat{\mathbf{K}}^{1:L}$ , we can show that the rate of requesting jobs is approximately equal to the arrival rate for each job type. Therefore, the number of tokens rarely reaches 0 or  $\eta_{\max}$ . This implies the rarity of virtual jobs and jobs on backup servers. The proofs are provided in Appendix B.2.

### 5.3 Proof of Theorem 3 Assuming Irreducibility Based on Lemmas 1–4.

**PROOF.** First we show that Lemmas 2 and 3 imply the closeness between  $\mathbf{K}^{1:L}$  and  $\overline{\mathbf{K}}^{1:L}$ . By Lemma 2, for any  $f \in \text{Lip}(1)$ , we have  $\mathbb{E}[f(\overline{\mathbf{K}}^{1:L})] - \mathbb{E}[f(\widehat{\mathbf{K}}^{1:L})] = O(\sqrt{r})$ . By Lemma 3,  $\mathbb{E} \left[ \sum_{\ell=1}^L \zeta_i^\ell \right] = O(\sqrt{r})$ . Recall that  $\widehat{\mathbf{K}}^\ell = \mathbf{K}^\ell + \boldsymbol{\zeta}^\ell$ , so  $\mathbb{E}[f(\widehat{\mathbf{K}}^{1:L})] - \mathbb{E}[f(\mathbf{K}^{1:L})] = O(\sqrt{r})$ . Therefore,

$$\mathbb{E} \left[ f(\overline{\mathbf{K}}^{1:L}) \right] - \mathbb{E} \left[ f(\mathbf{K}^{1:L}) \right] = O \left( \sqrt{r} \right). \quad (9)$$

Now we prove (5), the bound on the expected number of the active servers, by taking a suitable  $f$  in (9). Observe that

$$\sum_{\mathbf{k} \neq 0} \mathbb{E}[X_{\mathbf{k}}] - L \cdot \mathbb{P}(\overline{\mathbf{K}} \neq 0) = \mathbb{E} \left[ \sum_{\ell=1}^L \mathbb{1}_{\{\mathbf{K}^\ell \neq 0\}} \right] - \mathbb{E} \left[ \sum_{\ell=1}^L \mathbb{1}_{\{\overline{\mathbf{K}}^\ell \neq 0\}} \right] + \mathbb{E} \left[ \sum_{\ell=L+1}^\infty \mathbb{1}_{\{\mathbf{K}^\ell \neq 0\}} \right], \quad (10)$$

where the last term on RHS is  $O(\sqrt{r})$  by Lemma 4. To show that the difference between the first two terms on the RHS of (10) are also  $O(\sqrt{r})$ , consider  $f_1(\mathbf{k}^{1:L}) \triangleq \sum_{\ell=1}^L \mathbb{1}_{\{\mathbf{k}^\ell \neq 0\}}$ . Because

$$|f_1(\mathbf{k}^{1:L}) - f_1(\mathbf{k}'^{1:L})| = \left| \sum_{\ell=1}^L (\mathbb{1}_{\{\mathbf{k}^\ell \neq 0\}} - \mathbb{1}_{\{\mathbf{k}'^\ell \neq 0\}}) \right| \leq \sum_{\ell=1}^L \mathbb{1}_{\{\mathbf{k}^\ell \neq \mathbf{k}'^\ell\}} \leq \|\mathbf{k}^{1:L} - \mathbf{k}'^{1:L}\|,$$

for any  $\mathbf{k}^{1:L}, \mathbf{k}'^{1:L} \in \mathcal{K}^L$ , we have  $f_1 \in \text{Lip}(1)$ . By (9),  $\mathbb{E}[f_1(\mathbf{K}^\ell)] - \mathbb{E} \left[ \sum_{\ell=1}^L f_1(\overline{\mathbf{K}}^\ell) \right] = O(\sqrt{r})$ . Therefore,  $\sum_{\mathbf{k} \neq 0} \mathbb{E}[X_{\mathbf{k}}] - L \cdot \mathbb{P}(\overline{\mathbf{K}} \neq 0) = O(\sqrt{r})$ . Recall that  $L = \lceil \overline{N} \rceil$ , so we get (5).

Similarly, to prove (6), we observe that

$$\sum_{\mathbf{k} \neq 0} h(\mathbf{k}) \mathbb{E}[X_{\mathbf{k}}] - L \cdot \mathbb{E}[h(\overline{\mathbf{K}})] = \mathbb{E} \left[ \sum_{\ell=1}^L h(\mathbf{K}^\ell) \right] - \mathbb{E} \left[ \sum_{\ell=1}^L h(\overline{\mathbf{K}}^\ell) \right] + \mathbb{E} \left[ \sum_{\ell=L+1}^\infty h(\mathbf{K}^\ell) \right]. \quad (11)$$

The last term of (11) can be bounded as  $\mathbb{E}\left[\sum_{\ell=L+1}^{\infty} h(\mathbf{K}^{\ell})\right] \leq \mathbb{E}\left[\sum_{\ell=L+1}^{\infty} \mathbb{1}_{\{K_i^{\ell} \neq 0\}}\right] \cdot \max_{\mathbf{k} \in \mathcal{K}} h(\mathbf{k})$ , which is  $O(\sqrt{r})$  by Lemma 4 and the fact that  $\mathcal{K}$  is a finite set. To show that the difference between the first two terms on the RHS of (11) is also  $O(\sqrt{r})$ , consider  $f_2(\mathbf{k}^{1:L}) = \frac{1}{\Gamma} \sum_{\ell=1}^L h(\mathbf{k}^{\ell})$ , where  $\Gamma$  is the Lipschitz constant of  $h(\cdot)$ . Because

$$|f_2(\mathbf{k}^{1:L}) - f_2(\mathbf{k}'^{1:L})| = \frac{1}{\Gamma} \left| \sum_{\ell=1}^L (h(\mathbf{k}^{\ell}) - h(\mathbf{k}'^{\ell})) \right| \leq \sum_{\ell=1}^L \|\mathbf{k}^{\ell} - \mathbf{k}'^{\ell}\| = \|\mathbf{k}^{1:L} - \mathbf{k}'^{1:L}\|,$$

for any  $\mathbf{k}^{1:L}, \mathbf{k}'^{1:L} \in \mathcal{K}^L$ , we have  $f_2 \in \text{Lip}(1)$ . By (9),  $\mathbb{E}[f_2(\mathbf{K}^{\ell})] - \mathbb{E}\left[\sum_{\ell=1}^L f_2(\bar{\mathbf{K}}^{\ell})\right] = O(\sqrt{r})$ . Therefore,  $\sum_{\mathbf{k} \neq 0} h(\mathbf{k})\mathbb{E}[X_{\mathbf{k}}] - L \cdot \mathbb{E}[h(\bar{\mathbf{K}})] = O(\sqrt{r})$ . Recall that  $L = \lceil \bar{N} \rceil$ , so we get (6).

To show (7) and (8), noting that  $\lceil \bar{N} \rceil = \Theta(r)$ , we have

$$\begin{aligned} N(\sigma) &= \sum_{\mathbf{k} \neq 0} \mathbb{E}[X_{\mathbf{k}}] \leq \lceil \bar{N} \rceil + O(\sqrt{r}) = (1 + O(r^{-0.5})) \cdot \bar{N}, \\ C(\sigma) &= \frac{\sum_{\mathbf{k} \neq 0} h(\mathbf{k})\mathbb{E}[X_{\mathbf{k}}]}{\sum_{\mathbf{k} \neq 0} \mathbb{E}[X_{\mathbf{k}}]} = \frac{\sum_{\mathbf{k} \neq 0} h(\mathbf{k})\pi(\mathbf{k}) + O(r^{-0.5})}{1 - \pi(\mathbf{0}) + O(r^{-0.5})} \leq (1 + O(r^{-0.5})) \cdot \epsilon, \end{aligned}$$

where in the last inequality we have used the fact that  $\epsilon > 0$ . This completes the proof.  $\square$

#### 5.4 More Details on the System and Proof of Lemma 2

To bound the distance between  $\bar{\mathbf{K}}^{1:L}$  and  $\hat{\mathbf{K}}^{1:L}$ , observe that because  $f \in \text{Lip}(1)$  and  $\sum_{\ell=1}^L \sum_{i \in \mathcal{I}} \eta_i^{\ell} = O(\sqrt{r})$ , it suffices to bound the Wasserstein distance between  $\bar{\mathbf{K}}^{1:L}$  and  $\hat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}$ , i.e.,

$$\sup_{f \in \text{Lip}(1)} \left\{ \mathbb{E}\left[f(\bar{\mathbf{K}}^{1:L})\right] - \mathbb{E}\left[f(\hat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L})\right] \right\} = O(\sqrt{r}), \quad (12)$$

where  $f(\hat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L})$  is a valid expression because  $\hat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L} \in \mathcal{K}$  as discussed in Remark 1 below.

**5.4.1 More Details on System Dynamics and Generator.** To prepare for the proof, we first look into the dynamics of the two systems under study. In particular, we write out the *generators* of  $\bar{\mathbf{K}}^{1:L}(t)$  and  $(\hat{\mathbf{K}}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$ , which are used in the Stein's method arguments.

We first examine the dynamics of the single-server system under Markovian policy  $\bar{\sigma}$ . Four types of events change a single-server system's configuration: internal transitions, departures, reactive requests, and proactive requests (see Section 4.1). The change of configuration due to any event can be represented by the diagram

$$\mathbf{k} \rightarrow \mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a},$$

where the arrow  $\mathbf{k} \rightarrow \mathbf{k}'$  denotes an internal transition or a departure from configuration  $\mathbf{k}$  to  $\mathbf{k}'$  if  $\mathbf{k} \neq \mathbf{k}'$ ; the arrow  $\mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a}$  denotes a reactive request that adds  $\mathbf{a}$  jobs to the system if  $\mathbf{k} \neq \mathbf{k}'$ , and denotes a proactive request if  $\mathbf{k} = \mathbf{k}'$ . We call the above change of configuration a *transition*, and denote its rate as  $\gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$ . Let  $E(\mathbf{k})$  denote the set of possible  $(\mathbf{k}', \mathbf{a})$  pairs in a transition.

We define the total transition rate at configuration  $\mathbf{k}$  as  $\gamma_{\mathbf{k}} \triangleq \sum_{(\mathbf{k}',\mathbf{a}) \in E(\mathbf{k})} \gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$ , and define the maximal transition rate  $\gamma_{\max} = \max_{\mathbf{k} \in \mathcal{K}} \gamma_{\mathbf{k}}$ . Since  $\mathcal{K}$  is a finite set, we have  $\gamma_{\max} < \infty$ . Also, observe that the request rate of type  $i$  jobs is given by

$$\bar{\lambda}_i \triangleq \sum_{\mathbf{k}} \sum_{(\mathbf{k}',\mathbf{a}) \in E(\mathbf{k})} \gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})} a_i \cdot \pi(\mathbf{k}), \quad (13)$$

where  $\pi$  denotes the stationary distribution of single-server configuration under policy  $\bar{\sigma}$ .

Next, we focus on the dynamics of  $L$  i.i.d. copies of single-server systems. Consider the generator  $\bar{G}$  of the corresponding Markov chain  $\{\bar{\mathbf{K}}^{1:L}(t)\}$ , which is a linear operator on functions  $g: \mathcal{K}^L \rightarrow \mathbb{R}$  defined as:

$$\bar{G}g(\mathbf{k}^{1:L}) \triangleq \frac{d}{dt} \mathbb{E} \left[ g \left( \bar{\mathbf{K}}^{1:L}(t) \right) \middle| \bar{\mathbf{K}}^{1:L}(0) = \mathbf{k}^{1:L} \right] \Big|_{t=0}, \quad (14)$$

and we call the resulting function  $\bar{G}g(\cdot)$  the *drift* of  $g(\cdot)$ . Based on the transition rates defined above, we have

$$\bar{G}g(\mathbf{k}^{1:L}) = \sum_{\ell=1}^L \sum_{(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k}^\ell)} \gamma_{\mathbf{k}^\ell, (\mathbf{k}', \mathbf{a})} (g(\cdot, \mathbf{k}' + \mathbf{a}, \cdot) - g(\cdot, \mathbf{k}^\ell, \cdot)), \quad (15)$$

where  $g(\cdot, \mathbf{k}' + \mathbf{a}, \cdot) - g(\cdot, \mathbf{k}^\ell, \cdot)$  is a shorthand for  $g(\mathbf{k}^1, \dots, \mathbf{k}^{\ell-1}, \mathbf{k}' + \mathbf{a}, \mathbf{k}^{\ell+1}, \dots, \mathbf{k}^L) - g(\mathbf{k}^{1:L})$ , i.e., we use  $\cdot$  to omit the entries that agree with  $\mathbf{k}^{1:L}$ .

Similarly, for the infinite-server system, consider the generator  $\hat{G}$  of  $(\hat{\mathbf{K}}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  defined as

$$\hat{G}\psi(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L}) \triangleq \frac{d}{dt} \mathbb{E} \left[ \psi \left( \hat{\mathbf{K}}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t) \right) \middle| \hat{\mathbf{K}}^{1:L}(0) = \mathbf{k}^{1:L}, \mathbf{K}^{1:L}(0) = \boldsymbol{\eta}^{1:L} \right] \Big|_{t=0}, \quad (16)$$

for any function  $\psi: (\mathcal{K} \times \mathcal{K})^L \rightarrow \mathbb{R}$ . The drift of  $\psi$  under  $\hat{G}$  turns out to have a similar decoupled form as  $\bar{G}g$ : observe that for each  $\ell$ , the transition of  $(\hat{\mathbf{K}}^\ell(t), \boldsymbol{\eta}^\ell(t))$  from  $(\mathbf{k}, \boldsymbol{\eta})$  to  $(\mathbf{k}', \boldsymbol{\eta} + \mathbf{a} \mathbb{1}_{\{\boldsymbol{\eta}=0\}})$  occurs at the rate  $\gamma_{\mathbf{k}, (\mathbf{k}', \mathbf{a})}$  for each  $(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k})$ , and any real or virtual job arrivals do not change the sum  $\hat{\mathbf{K}}^\ell(t) + \boldsymbol{\eta}^\ell(t)$ . Consider any function  $g: \mathcal{K}^L \rightarrow \mathbb{R}$  and the function  $\psi(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L}) = g(\mathbf{k}^{1:L} + \boldsymbol{\eta}^{1:L})$ .

$$\begin{aligned} \hat{G}\psi(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L}) &= \sum_{\ell=1}^L \sum_{(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k}^\ell)} \gamma_{\mathbf{k}^\ell, (\mathbf{k}', \mathbf{a})} (g(\cdot, \mathbf{k}' + \mathbf{a}, \cdot) - g(\cdot, \mathbf{k}^\ell, \cdot)) \mathbb{1}_{\{\boldsymbol{\eta}^\ell=0\}} \\ &\quad + \sum_{\ell=1}^L \sum_{(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k}^\ell)} \gamma_{\mathbf{k}^\ell, (\mathbf{k}', \mathbf{a})} (g(\cdot, \mathbf{k}' + \boldsymbol{\eta}^\ell, \cdot) - g(\cdot, \mathbf{k}^\ell + \boldsymbol{\eta}^\ell, \cdot)) \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq 0\}}. \end{aligned} \quad (17)$$

In this context  $g(\cdot, \mathbf{k}' + \boldsymbol{\eta}^\ell, \cdot) - g(\cdot, \mathbf{k}^\ell + \boldsymbol{\eta}^\ell, \cdot)$  is a shorthand for  $g(\mathbf{k}^1 + \boldsymbol{\eta}^1, \dots, \mathbf{k}^{\ell-1} + \boldsymbol{\eta}^{\ell-1}, \mathbf{k}' + \boldsymbol{\eta}^\ell, \mathbf{k}^{\ell+1} + \boldsymbol{\eta}^{\ell+1}, \dots, \mathbf{k}^L + \boldsymbol{\eta}^L) - g(\mathbf{k}^{1:L} + \boldsymbol{\eta}^{1:L})$ . In other words, we use  $\cdot$  to omit the entries of  $g$ 's input that agree with the corresponding entries of  $\mathbf{k}^{1:L} + \boldsymbol{\eta}^{1:L}$ .

*Remark 1.* In (17), although  $g$  is only defined on the domain  $\mathcal{K}^L$ , it is valid to write  $\mathbf{k}^{1:L} + \boldsymbol{\eta}^{1:L}$  as its input because we always have  $\hat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell \in \mathcal{K}$ , i.e., the total number of real jobs, virtual jobs, and tokens on a normal server never exceeds  $K_{\max}$ . To see why this is true, the single-server policy  $\bar{\sigma}$  requests jobs only when there are no tokens on the server, and it will not request more than  $K_{\max} - n$  jobs if there are already  $n$  real and virtual jobs on the server.

#### 5.4.2 Proof of Lemma 2.

**PROOF. Generator Comparison.** For any  $f \in \text{Lip}(1)$ , consider the Poisson equation (see, e.g., [3]) that solves for  $g_f: \mathcal{K}^L \rightarrow \mathbb{R}$ :

$$\mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L}) \right] - f(\mathbf{k}^{1:L}) = \bar{G}g_f(\mathbf{k}^{1:L}). \quad (18)$$

We let  $\mathbf{k}^{1:L} = \hat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}$  in (18) and take the expectation. This results in

$$\mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L}) \right] - \mathbb{E} \left[ f(\hat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) \right] = \mathbb{E} \left[ \bar{G}g_f(\hat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) \right]. \quad (19)$$

On the other hand, because  $(\hat{\mathbf{K}}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  is a finite-state Markov chain, we have

$$\mathbb{E} \left[ \hat{G}\psi_f(\hat{\mathbf{K}}^{1:L}, \boldsymbol{\eta}^{1:L}) \right] = 0, \quad (20)$$

where  $\psi_f$  is given by  $\psi_f(\widehat{\mathbf{K}}^{1:L}, \boldsymbol{\eta}^{1:L}) = g_f(\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L})$ . Subtracting (20) from (19), we get

$$\mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L}) \right] - \mathbb{E} \left[ f(\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) \right] = \mathbb{E} \left[ \bar{G}g_f(\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) - \widehat{G}\psi_f(\widehat{\mathbf{K}}^{1:L}, \boldsymbol{\eta}^{1:L}) \right]. \quad (21)$$

We want to show that  $\bar{G}$  and  $\widehat{G}$  are close so that we can bound the RHS of (21).

Now we plug the formula of the generators in (15) and (17) into the RHS of (21) and get

$$\begin{aligned} & \left| \bar{G}g_f(\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) - \widehat{G}\psi_f(\widehat{\mathbf{K}}^{1:L}, \boldsymbol{\eta}^{1:L}) \right| \\ & \stackrel{(i)}{=} \left| \sum_{\ell=1}^L \sum_{(\mathbf{k}', \mathbf{a}) \in E(\widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell)} \gamma_{\widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell, (\mathbf{k}', \mathbf{a})} \left( g_f(\cdot, \mathbf{k}' + \mathbf{a}, \cdot) - g_f(\cdot, \widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell, \cdot) \right) \cdot \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq \mathbf{0}\}} \right. \\ & \quad \left. - \sum_{\ell=1}^L \sum_{(\mathbf{k}', \mathbf{a}) \in E(\widehat{\mathbf{K}}^\ell)} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} \left( g_f(\cdot, \mathbf{k}' + \boldsymbol{\eta}^\ell, \cdot) - g_f(\cdot, \widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell, \cdot) \right) \cdot \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq \mathbf{0}\}} \right| \\ & \stackrel{(ii)}{\leq} \sum_{\ell=1}^L \gamma_{\widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell} \cdot \sup_{(\mathbf{k}', \mathbf{a}) \in E(\widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell)} \left| g_f(\cdot, \mathbf{k}' + \mathbf{a}, \cdot) - g_f(\cdot, \widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell, \cdot) \right| \cdot \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq \mathbf{0}\}} \\ & \quad + \sum_{\ell=1}^L \gamma_{\widehat{\mathbf{K}}^\ell} \cdot \sup_{(\mathbf{k}', \mathbf{a}) \in E(\widehat{\mathbf{K}}^\ell)} \left| g_f(\cdot, \mathbf{k}' + \boldsymbol{\eta}^\ell, \cdot) - g_f(\cdot, \widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell, \cdot) \right| \cdot \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq \mathbf{0}\}} \\ & \stackrel{(iii)}{\leq} 2\gamma_{\max} \cdot \sum_{\ell=1}^L \sup_{\mathbf{k}' \in \mathcal{K}} \left| g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \widehat{\mathbf{K}}^\ell + \boldsymbol{\eta}^\ell, \cdot) \right| \cdot \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq \mathbf{0}\}} \\ & \leq 2\gamma_{\max} \cdot \sup_{\mathbf{k}, \mathbf{k}' \in \mathcal{K}} \left| g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \mathbf{k}, \cdot) \right| \cdot \sum_{\ell=1}^L \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq \mathbf{0}\}}, \end{aligned} \quad (22)$$

where in  $g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \mathbf{k}, \cdot)$  we have omitted the entries that agree with  $\mathbf{k}^{1:L}$ . The equality (i) is true because each of the  $\ell$ -th terms in  $\bar{G}$  and  $\widehat{G}$  is equal if  $\boldsymbol{\eta}^\ell = \mathbf{0}$ . For the inequalities (ii) and (iii), recall that  $\gamma_k$  is the total transition rate given by  $\gamma_k \triangleq \sum_{(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k})} \gamma_{\mathbf{k}, (\mathbf{k}', \mathbf{a})}$ , and  $\gamma_{\max} = \max_{\mathbf{k} \in \mathcal{K}} \gamma_k$ . Observe that

$$\sum_{\ell=1}^L \mathbb{1}_{\{\boldsymbol{\eta}^\ell \neq \mathbf{0}\}} \leq \sum_{\ell=1}^L \sum_{i \in I} \eta_i^\ell \leq |I| \cdot \eta_{\max} = O(\sqrt{r}).$$

Therefore (22) can be further bounded by

$$\begin{aligned} \left| \bar{G}g_f(\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) - \widehat{G}\psi_f(\widehat{\mathbf{K}}^{1:L}, \boldsymbol{\eta}^{1:L}) \right| & \leq 2\gamma_{\max} \cdot \sup_{\mathbf{k}, \mathbf{k}' \in \mathcal{K}} \left| g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \mathbf{k}, \cdot) \right| \cdot \sum_{\ell=1}^L \sum_{i \in I} \eta_i^\ell \\ & \leq 2\gamma_{\max} \cdot \sup_{\mathbf{k}, \mathbf{k}' \in \mathcal{K}} \left| g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \mathbf{k}, \cdot) \right| \cdot O(\sqrt{r}). \end{aligned}$$

To prove (12), it remains to show that

$$\sup_{\mathbf{k}, \mathbf{k}' \in \mathcal{K}} \left| g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \mathbf{k}, \cdot) \right| = O(1). \quad (23)$$

**Stein Factor Bound.** To prove (23), observe that the following  $g_f(\cdot)$  is a solution to the Poisson equation (18):

$$g_f(\mathbf{k}^{1:L}) = \int_0^\infty \mathbb{E} \left[ \left( f(\bar{\mathbf{K}}^{1:L}(t)) - \mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L}) \right] \right) \middle| \bar{\mathbf{K}}^{1:L}(0) = \mathbf{k}^{1:L} \right] dt. \quad (24)$$

This allows us to bound the difference of  $g_f$  using coupling. Specifically, we define the coupling of two systems, each consisting of  $L$  i.i.d. copies of the single-server system under  $\bar{\sigma}$ . The two systems are initialized with configurations  $(\cdot, \mathbf{k}', \cdot)$  and  $(\cdot, \mathbf{k}, \cdot)$  that only differ at the  $\ell$ -th server, where we omit the entries that agree with  $\mathbf{k}^{1:L}$ . Let  $(\bar{\mathbf{K}}^{1:L,1}(t), \bar{\mathbf{K}}^{1:L,2}(t))$  be the joint configuration of the two systems, which is actually  $2L$  i.i.d. copies of the single-server system. As a result, we can specify the couplings  $(\bar{\mathbf{K}}^{\ell',1}(t), \bar{\mathbf{K}}^{\ell',2}(t))$  for different  $\ell'$  separately. For  $\ell' \neq \ell$ , the corresponding server in the two systems have the same initial configurations, so we can always keep their configurations identical. For the  $\ell$ -th servers, we let them evolve independently following their own dynamics until a stopping time  $\tau_{\text{mix}}$  when their configurations become the same. After that, we can use coupling to keep their configurations identical. Under this coupling, it is not hard to see that

$$\begin{aligned} |g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \mathbf{k}, \cdot)| &= \left| \int_0^\infty \mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L,1}(t)) - f(\bar{\mathbf{K}}^{1:L,2}(t)) \right] dt \right| \\ &\leq \mathbb{E} \left[ \int_0^\infty \left| f(\bar{\mathbf{K}}^{1:L,1}(t)) - f(\bar{\mathbf{K}}^{1:L,2}(t)) \right| dt \right] \\ &\leq \mathbb{E} \left[ \int_0^\infty \sum_{\ell'=1}^L \|\bar{\mathbf{K}}^{\ell',1}(t) - \bar{\mathbf{K}}^{\ell',2}(t)\| dt \right] \\ &= \mathbb{E} \left[ \int_0^{\tau_{\text{mix}}} \|\bar{\mathbf{K}}^{\ell,1}(t) - \bar{\mathbf{K}}^{\ell,2}(t)\| dt \right], \end{aligned} \quad (25)$$

where in the second inequality we have used the fact that  $f$  is 1-Lipschitz continuous under the  $L^1$  norm of the space  $\mathcal{K}^L$ . For each pair of  $\mathbf{k}, \mathbf{k}'$ , observe that because  $\bar{\sigma}$  is a  $\mathbf{k}^0$ -irreducible policy,  $\mathbb{E}[\tau_{\text{mix}}]$  is finite; and because  $\mathcal{K}$  is a finite set,  $\|\bar{\mathbf{K}}^{\ell,1}(t) - \bar{\mathbf{K}}^{\ell,2}(t)\|$  is uniformly bounded. All these finite quantities depend on a single-server system under a policy  $\bar{\sigma}$  that is independent of  $r$ . As a result, the last expression in (25) is of constant order. Moreover, because there are finite pairs of  $(\mathbf{k}, \mathbf{k}')$ , the supremum  $\sup_{\mathbf{k}, \mathbf{k}'} \mathbb{E} \left[ \int_0^{\tau_{\text{mix}}} \|\bar{\mathbf{K}}^{\ell,1}(t) - \bar{\mathbf{K}}^{\ell,2}(t)\| dt \right]$  is also of constant order, independent of  $r$ . This proves the Stein factor bound in (23). Together with the generator comparison, we have proved

$$\sup_{f \in \text{Lip}(1)} \mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L}) \right] - \mathbb{E} \left[ f(\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) \right] = O(\sqrt{r}). \quad (12)$$

Because  $\sum_{\ell=1}^L \sum_{i \in I} \eta_i^\ell \leq |I| \eta_{\max} = O(\sqrt{r})$ , for any  $f \in \text{Lip}(1)$ , we have

$$\left| \mathbb{E} \left[ f(\widehat{\mathbf{K}}^{1:L}) \right] - \mathbb{E} \left[ f(\widehat{\mathbf{K}}^{1:L} + \boldsymbol{\eta}^{1:L}) \right] \right| \leq \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{i \in I} \eta_i^\ell(t) \right] = O(\sqrt{r}). \quad (26)$$

Plugging the above equation to (12), we get  $\sup_{f \in \text{Lip}(1)} \mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L}) \right] - \mathbb{E} \left[ f(\widehat{\mathbf{K}}^{1:L}) \right] = O(\sqrt{r})$ . This proves Lemma 2.  $\square$

## 5.5 Role of Tokens and Virtual Jobs

This section aims to shed light on the role of tokens and virtual jobs in the proposed policy, JRS. We first outline how we devise the token-and-virtual-job mechanism from the perspective of generators. To begin with, consider the scenario where dispatch decisions are solely based on real-job configurations  $(\mathbf{K}^\ell)_{\ell=1,2,\dots}$ . In this case, the transitions of servers' configurations would be correlated in general due to job arrivals, which are dispatched based on the joint configuration of all servers. To break this correlation, let us introduce tokens but not set an upper limit yet on the number of tokens (thus no virtual jobs). Observe that  $\mathbf{K}^{1:L}(t) + \boldsymbol{\eta}^{1:L}(t)$  remains unchanged by



job arrivals, so tokens remove the correlation brought about by job arrivals. However, because tokens lack internal phase transitions or departures, the transition dynamics of  $K^{1:L}(t) + \eta^{1:L}(t)$  will diverge from  $\bar{K}^{1:L}(t)$  when  $\eta(t)$  is large. In other words, although tokens help decouple the transitions on servers, they cannot keep the transitions of  $K^{1:L}(t) + \eta^{1:L}(t)$  close to  $\bar{K}^{1:L}(t)$ . To solve this issue, we finally introduce the mechanism of converting tokens into virtual jobs when the number of tokens is high, where virtual jobs can make internal phase transitions or departures just like real jobs. Now, the sum  $K^{1:L}(t) + \zeta^{1:L}(t) + \eta^{1:L}(t)$  remains unchanged by job arrivals nor the creation of virtual jobs, and the internal phase transitions and job departures are similar to those of  $\bar{K}^{1:L}(t)$ . More formally, the generators of  $K^{1:L}(t) + \zeta^{1:L}(t) + \eta^{1:L}(t)$  and  $\bar{K}^{1:L}(t)$  are close to each other – their additive difference can be upper bounded by a quantity proportional to the expected number of tokens, as shown in (22). Therefore, by regulating the number of tokens, we can control the difference between the generators of  $K^{1:L}(t) + \zeta^{1:L}(t) + \eta^{1:L}(t)$  and  $\bar{K}^{1:L}(t)$ .

Another key design component of JRS is that the subroutine requests jobs based on the observed configurations. This has been used in the proof of Lemma 2 to show that the observed configurations  $K^{1:L} + \zeta^{1:L}$  are close to  $\bar{K}^{1:L}$ , which consist of  $L$  i.i.d. single-server systems. Recall that each single-server system in  $\bar{K}^{1:L}$  is designed to have a throughput of  $\lambda_i r / L$  for each job type  $i \in \mathcal{I}$ . Therefore, the proximity between  $K^{1:L} + \zeta^{1:L}$  and  $\bar{K}^{1:L}$  ensures that the job request rate mirrors the arrival rate for each job type, regardless of the real-job configurations. The fact that these two rates are approximately equal is important for proving Lemma 3 and Lemma 4. It guarantees that both the rate of generating virtual jobs (when there are too many tokens) and the rate of dispatching jobs to backup servers (when there are no tokens) are appropriate.

A natural follow-up question is whether the usage of tokens and virtual jobs is fundamental or an artifact of our analysis technique. For example, it is unclear whether removing the upper limit on the number of tokens would still yield an asymptotically optimal policy. This is an interesting question that we do not have a complete answer to. The token-and-virtual-job mechanism emerges as a natural choice under our analysis framework. Nevertheless, it is worth noting that our analysis primarily treats each server's configuration as a generic Markov chain, without utilizing many properties specific to the stochastic bin-packing setting. An exception to this is the proofs in Appendix B.2, where we use the model that each job leaves the system within a constant expected time. It would be interesting to explore more properties of the problem to better understand policy designs without auxiliary state variables like tokens and virtual jobs.

## 6 CONCLUSION

In this paper, we study a new setting of stochastic bin-packing in service systems that features time-varying item sizes. Since our formulation is motivated by the problem of virtual-machine scheduling in computing systems, we use the terminology of jobs and servers, where jobs are viewed as items, whose sizes are their resource requirements, and servers as bins. The time-varying item sizes capture the emerging trend in practice that jobs' resource requirements vary over time. Our goal is to design a job dispatch policy to minimize the expected number of active servers in steady state, subject to a constraint on resource contentions. Our main result is the design of a policy that achieves an optimality gap of  $O(\sqrt{r})$ , where  $r$  is the scaling factor of the arrival rate. When specialized to the setting where jobs' resource requirements remain fixed over time, this result improves upon the state-of-the-art  $o(r)$  optimality gap. Our technical approach highlights a novel policy conversion framework, JOIN-REQUESTING-SERVER, that reduces the policy design problem to that in a single-server system.

There are several potential directions that may be worth further exploration. One direction is to strengthen the optimality result within the current setting. Specifically, it is interesting to investigate: (i) whether it is possible to achieve an optimality gap smaller than  $\Theta(\sqrt{r})$ ; and (ii) whether there exist asymptotically optimal policies whose average cost rate of resource contention satisfies the budget strictly instead of asymptotically.

We are also interested in extending our technique to the optimal control of other systems with similar structures. Intuitively, this technique could be applied to systems with many components that evolve mostly independently but are weakly coupled by certain constraints. Viewing each component as a server, we can define a suitable single-server problem and then design a policy for the original system to track the dynamics of the optimal single-server solution. Below we list several variations of our model that can potentially be analyzed using the proposed technique.

- A model where jobs running on each server will be put into a local queue when there are resource contentions. The goal thus becomes finding the optimal trade-offs between the number of active servers and the waiting time of the jobs.
- A model that allows each server to have a Markovian state that affects the dynamics of the jobs running on the server.
- A model that allows jobs to migrate to different servers at the cost of migration delays.
- A closed-system model where jobs re-enter the system after completion.

A third possible direction is to tackle the problem when the arrival rates and the parameters in the job model are unknown, as mentioned in Section 4.3. A possible approach is to develop an approximate version of the JRS framework, where the optimal single-server policy and the simulator for the virtual jobs are both learned from data. It is desirable to design such an approximate framework whose performance degrades gracefully as the approximation error increases.

## ACKNOWLEDGMENTS

Y. Hong and W. Wang are supported in part by NSF grants CNS-200773 and ECCS-2145713. Q. Xie is supported in part by NSF grant CNS-1955997.

## REFERENCES

- [1] Nikhil Ayyadevara, Rajni Dabas, Arindam Khan, and K. V. N. Sreenivas. 2022. Near-Optimal Algorithms for Stochastic Online Bin Packing. In *Proc. Int. Conf. Automata, Languages and Programming (ICALP)*, Vol. 229. 12:1–12:20.
- [2] Noman Bashir, Nan Deng, Krzysztof Rządca, David Irwin, Sree Kodak, and Rohit Jnagal. 2021. Take It to the Limit: Peak Prediction-Driven Resource Overcommitment in Datacenters. In *Proc. European Conf. Computer Systems (EuroSys)*. Online Event, United Kingdom, 556–573.
- [3] Anton Braverman. 2022. The Prelimit Generator Comparison Approach of Stein’s Method. *Stoch. Syst.* 12, 2 (2022), 181–204.
- [4] Anton Braverman and J. G. Dai. 2017. Stein’s method for steady-state diffusion approximations of  $M/Ph/n + M$  systems. *Ann. Appl. Probab.* 27 (Feb. 2017), 550–581.
- [5] Anton Braverman, J. G. Dai, and Jiekun Feng. 2017. Stein’s method for steady-state diffusion approximations: an introduction through the Erlang-A and Erlang-C models. *Stoch. Syst.* 6, 2 (2017), 301–366.
- [6] Niv Buchbinder, Yaron Fairstein, Konstantina Mellou, Ishai Menache, and Joseph (Seffi) Naor. 2021. Online Virtual Machine Allocation with Lifetime and Load Predictions. *ACM SIGMETRICS Perform. Eval. Rev.* 49, 1 (May 2021), 9–10.
- [7] Google Cloud. 2023. Overcommitting CPUs on sole-tenant VMs. <https://cloud.google.com/compute/docs/nodes/overcommitting-cpus-sole-tenant-vm>.
- [8] Google Cloud. 2023. Virtual machine instances. <https://cloud.google.com/compute/docs/instances>.
- [9] E. G. Coffman, Jr., M. R. Garey, and D. S. Johnson. 1983. Dynamic Bin Packing. *SIAM J. Comput.* 12, 2 (1983), 227–258.
- [10] Coastas Courcoubetis and Richard Weber. 1990. Stability of On-Line Bin Packing with Random Arrivals and Long-Run-Average Constraints. *Probab. Eng. Inf. Sci.* 4, 4 (1990), 447–460.
- [11] C. Courcoubetis and R. R. Weber. 1986. Necessary and Sufficient Conditions for Stability of a Bin-Packing System. *J. Appl. Probab.* 23, 4 (1986), 989–999.

- [12] Janos Csirik, David S. Johnson, Claire Kenyon, James B. Orlin, Peter W. Shor, and Richard R. Weber. 2006. On the Sum-of-Squares Algorithm for Bin Packing. *J. ACM* 53, 1 (Jan. 2006), 1–65.
- [13] Christina Delimitrou and Christos Kozyrakis. 2014. Quasar: Resource-Efficient and QoS-Aware Cluster Management. In *Proc. Int. Conf. Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. Salt Lake City, UT, 127–144.
- [14] Apache Software Foundation. 2023. Apache Mesos: Containerizers. <https://mesos.apache.org/documentation/latest/containerizers/>.
- [15] Apache Software Foundation. 2023. Apache Mesos: Oversubscription. <https://mesos.apache.org/documentation/latest/oversubscription/>.
- [16] Daniel Freund and Siddhartha Banerjee. 2019. Good prophets know when the end is near. *Available at SSRN: https://ssrn.com/abstract=3479189* (Nov. 2019).
- [17] Javad Ghaderi, Yuan Zhong, and R. Srikant. 2014. Asymptotic Optimality of BestFit for Stochastic Bin Packing. *SIGMETRICS Perform. Eval. Rev.* 42, 2 (Sept. 2014), 64–66.
- [18] Varun Gupta and Ana Radovanović. 2020. Interior-Point-Based Online Stochastic Bin Packing. *Oper. Res.* 68, 5 (2020), 1474–1492.
- [19] Leonard Kleinrock. 1975. *Queueing Systems*. John Wiley & Son.
- [20] Yusen Li, Xueyan Tang, and Wentong Cai. 2014. On Dynamic Bin Packing for Resource Allocation in the Cloud. In *Proc. Ann. ACM Symp. Parallelism in Algorithms and Architectures (SPAA)*. Prague, Czech Republic, 2–11.
- [21] David Lo, Liquan Cheng, Rama Govindaraju, Parthasarathy Ranganathan, and Christos Kozyrakis. 2015. Heracles: Improving resource efficiency at scale. In *Proc. ACM/IEEE Ann. Int. Symp. Computer Architecture (ISCA)*. Portland, OR, 450–462.
- [22] Siva Theja Maguluri and R. Srikant. 2013. Scheduling jobs with unknown duration in clouds. In *Proc. IEEE Int. Conf. Computer Communications (INFOCOM)*. Turin, Italy, 1887–1895.
- [23] Siva Theja Maguluri, R Srikant, and Lei Ying. 2012. Stochastic Models of Load Balancing and Scheduling in Cloud Computing Clusters. In *Proc. IEEE Int. Conf. Computer Communications (INFOCOM)*. Orlando, FL, 702–710.
- [24] Siva Theja Maguluri, R. Srikant, and Lei Ying. 2014. Heavy traffic optimal resource allocation algorithms for cloud computing clusters. *Perform. Eval.* 81 (2014), 20–39.
- [25] Sean Meyn. 2007. *Control Techniques for Complex Networks* (1st ed.). Cambridge University Press, USA.
- [26] Konstantinos Psychas and Javad Ghaderi. 2018. On Non-Preemptive VM Scheduling in the Cloud. In *Proc. ACM SIGMETRICS Int. Conf. Measurement and Modeling of Computer Systems*. Irvine, CA, 67–69.
- [27] Konstantinos Psychas and Javad Ghaderi. 2019. Scheduling Jobs with Random Resource Requirements in Computing Clusters. In *Proc. IEEE Int. Conf. Computer Communications (INFOCOM)*. 2269–2277.
- [28] Konstantinos Psychas and Javad Ghaderi. 2021. High-Throughput Bin Packing: Scheduling Jobs With Random Resource Demands in Clusters. *IEEE/ACM Trans. Netw.* 29, 1 (2021), 220–233.
- [29] Konstantinos Psychas and Javad Ghaderi. 2022. A Theory of Auto-Scaling for Resource Reservation in Cloud Services. *Stoch. Syst.* 12, 3 (2022), 227–252.
- [30] Charles Reiss, Alexey Tumanov, Gregory R. Ganger, Randy H. Katz, and Michael A. Kozuch. 2012. Heterogeneity and Dynamicity of Clouds at Scale: Google Trace Analysis. In *Proc. ACM Symp. Cloud Computing (SoCC)*. San Jose, CA, Article 7, 13 pages.
- [31] Krzysztof Rządca, Paweł Findeisen, Jacek Swiderski, Przemysław Zych, Przemysław Broniek, Jarek Kusmierek, Paweł Nowak, Beata Strack, Piotr Witusowski, Steven Hand, and John Wilkes. 2020. Autopilot: Workload Autoscaling at Google. In *Proc. European Conf. Computer Systems (EuroSys)*. Heraklion, Greece, Article 16, 16 pages.
- [32] Alexander L. Stolyar. 2013. An Infinite Server System with General Packing Constraints. *Oper. Res.* 61, 5 (2013), 1200–1217.
- [33] Alexander L. Stolyar. 2017. Large-scale heterogeneous service systems with general packing constraints. *Adv. Appl. Probab.* 49 (March 2017), 61–83. Issue 1.
- [34] Alexander L. Stolyar and Yuan Zhong. 2013. A large-scale service system with packing constraints: Minimizing the number of occupied servers. *ACM SIGMETRICS Perform. Eval. Rev.* 41, 1 (June 2013), 41–52.
- [35] Alexander L. Stolyar and Yuan Zhong. 2015. Asymptotic optimality of a greedy randomized algorithm in a large-scale service system with general packing constraints. *Queueing Syst.* 79 (June 2015), 117–143. Issue 2.
- [36] Alexander L. Stolyar and Yuan Zhong. 2021. A Service System with Packing Constraints: Greedy Randomized Algorithm Achieving Sublinear in Scale Optimality Gap. *Stoch. Syst.* 11 (June 2021), 83–111. Issue 2.
- [37] Muhammad Tirmazi, Adam Barker, Nan Deng, Md E. Haque, Zhijiang Gene Qin, Steven Hand, Mor Harchol-Balter, and John Wilkes. 2020. Borg: The next Generation. In *Proc. European Conf. Computer Systems (EuroSys)*. Heraklion, Greece, Article 30, 14 pages.
- [38] John Wilkes. 2019. Google cluster-usage traces v3. <http://github.com/google/cluster-data>.

[39] Qiaomin Xie, Xiaobo Dong, Yi Lu, and R. Srikant. 2015. Power of d Choices for Large-Scale Bin Packing: A Loss Model. In *Proc. ACM SIGMETRICS Int. Conf. Measurement and Modeling of Computer Systems*. Portland, OR, 321–334.

## A PROOF OF THEOREM 2 (LOWER BOUND)

PROOF. It is sufficient to show that given an infinite-server policy  $\sigma$  for  $\mathcal{P}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$  in (2), we have  $N(\sigma) \geq \bar{N}^*$ . To this end, we will construct a single-server policy  $\bar{\sigma}$  such that the resulting system configuration  $\bar{K}(\infty)$  in steady state satisfies:

$$\mathbb{E} [h(\bar{K}(\infty)) | \bar{K}(\infty) \neq \mathbf{0}] = C(\sigma) \leq \epsilon, \quad (27)$$

$$\bar{\lambda}_i = \frac{\lambda_i r}{N(\sigma)}, \quad \forall i \in \mathcal{I}. \quad (28)$$

Let  $\pi$  be the distribution of  $\bar{K}(\infty)$ , then  $(N(\sigma), \bar{\sigma}, \pi)$  is a feasible solution to the problem  $\bar{\mathcal{P}}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$  in (3). As a result, we have  $N(\sigma) \geq \bar{N}^*$ . Note that although  $\bar{\sigma}$  is actually non-Markovian, i.e., it makes decisions based on not only the current configuration but also the history, as we will show in Appendix C,  $\bar{N}^*$  is still a lower bound to the objective value that  $\bar{\sigma}$  can achieve in  $\bar{\mathcal{P}}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$ .

The construction of the single-server policy  $\bar{\sigma}$  involves simulating an infinite-server system under  $\sigma$  from the empty configuration. At time 0, the policy  $\bar{\sigma}$  randomly chooses the  $\ell$ -th server in the infinite-server system with probability  $p^\ell$ , for  $\ell = 1, 2, \dots$ . It then requests jobs for the single-server system according to a policy  $\bar{\sigma}^\ell$ . The key to our policy  $\bar{\sigma}^\ell$  is to make the single-server system emulate the job assignment at the  $\ell$ -th server of the simulated infinite-server system, but without incurring idleness. We first construct the policy  $\bar{\sigma}^\ell$ , and then specify the probabilities  $p^\ell$ .

Let us start by introducing some useful notation. Let  $\bar{K}^\ell(t)$  be the single-server system configuration under  $\bar{\sigma}^\ell$  at time  $t$  and  $K^\ell(t)$  be the configuration of the  $\ell$ -th simulated server in the infinite-server system under  $\sigma$ . We define a stochastic process  $\{s^\ell(t), t \geq 0\}$  as follows:  $s^\ell(t) = \max_\tau \{ \tau : \int_0^\tau \mathbb{1}_{\{K^\ell(x) \neq \mathbf{0}\}} dx = t \}$ . The “max” is well-defined because the integral is continuous in  $\tau$ . Intuitively,  $s^\ell(t)$  gives the maximum time when the accumulative busy time of the  $\ell$ -th server is  $t$ . Note that  $\{s^\ell(t), t \geq 0\}$  is only discontinuous when  $K^\ell(\tau)$  reaches 0, thus it is right-differentiable with derivative equal to 1 at any point.

We construct  $\bar{\sigma}^\ell$  and the simulation of the infinite-server system under  $\sigma$  in a way such that:

$$\bar{K}^\ell(t) = K^\ell(s^\ell(t)) \quad \forall t. \quad (29)$$

That is, we want that the single-server system has the same dynamic of the simulated  $\ell$ -th server except skipping the idle period. To this end, we couple the two systems as follows:

- (1) When the  $\ell$ -th simulated server  $K^\ell(s^\ell(t))$  receives a type  $i$  job, we let the single-server system  $\bar{K}^\ell(t)$  request a type  $i$  job at time  $t$ . For each such job, its phase transition process in the  $\ell$ -th simulated server is the same as that in the single-server system. That is, when we observe any internal transition or departure event in  $\bar{K}^\ell(t)$ , we produce a same event on the  $\ell$ -th simulated server  $K^\ell(s^\ell(t))$ .
- (2) The simulations of the rest of the infinite-server system under policy  $\sigma$  are driven by independently generated random seeds.

It is not hard to see that the simulated infinite server-system has the same stochastic behavior as an uncoupled system under  $\sigma$ . Moreover, as we couple all the events that happen in  $\bar{K}^\ell(t)$  and  $K^\ell(s^\ell(t))$ , together with the facts that  $\bar{K}^\ell(t)$  and  $K^\ell(s^\ell(t))$  are piecewise constant and  $\bar{K}^\ell(0-) = K^\ell(s^\ell(0-)) = \mathbf{0}$ , we get (29).

Next we claim that (29) implies the following relationship between the steady-state cost of the single-server system under  $\bar{\sigma}^\ell$  and the steady-state cost of the  $\ell$ -th simulated server in the

infinite-server system under  $\sigma$ :

$$\mathbb{E} \left[ h \left( \bar{K}^\ell(\infty) \right) \right] = \frac{\mathbb{E} \left[ h(K^\ell(\infty)) \right]}{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})}. \quad (30)$$

This is because for all  $\mathbf{k} \neq \mathbf{0}$ , we have

$$\begin{aligned} \frac{\mathbb{P}(K^\ell(\infty) = \mathbf{k})}{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})} &\stackrel{(a)}{=} \lim_{S \rightarrow \infty} \frac{\int_0^S \mathbb{1}_{\{K^\ell(s) = \mathbf{k}\}} ds}{\int_0^S \mathbb{1}_{\{K^\ell(s) \neq \mathbf{0}\}} ds} \stackrel{(b)}{=} \lim_{T \rightarrow \infty} \frac{\int_0^T \mathbb{1}_{\{K^\ell(s^\ell(t)) = \mathbf{k}\}} dt}{\int_0^T \mathbb{1}_{\{K^\ell(s^\ell(t)) \neq \mathbf{0}\}} dt} \\ &\stackrel{(c)}{=} \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \mathbb{1}_{\{K^\ell(s^\ell(t)) = \mathbf{k}\}} dt \stackrel{(d)}{=} \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \mathbb{1}_{\{\bar{K}^\ell(t) = \mathbf{k}\}} dt \\ &\stackrel{(e)}{=} \mathbb{P}(\bar{K}^\ell(\infty) = \mathbf{k}), \end{aligned}$$

where (a) and (e) hold because long-run averages converge to steady-state expectations; (b) is due to the fact that  $\int_0^T \mathbb{1}_{\{K^\ell(s^\ell(t)) = \mathbf{k}\}} dt = \int_0^{s^\ell(T)} \mathbb{1}_{\{K^\ell(s) = \mathbf{k}\}} ds$ , for any  $\mathbf{k} \neq \mathbf{0}$ ; (c) is due to the fact that  $\mathbb{1}_{\{K^\ell(s^\ell(t)) \neq \mathbf{0}\}} = 1$ ; and (d) follows from (29).

Let  $\bar{\lambda}_i^\ell$  be the long-run request rates of type  $i$  jobs in the single-server system under  $\bar{\sigma}^\ell$ , and  $\lambda_i^\ell$  be the throughput of type  $i$  jobs of  $\ell$ -th simulated server under  $\sigma$ . By the construction of  $\bar{\sigma}^\ell$ , the single-server system requests jobs based on the arrival events of the  $\ell$ -th simulated server, we have  $\bar{\lambda}_i^\ell = \frac{\lambda_i^\ell}{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})}$ ,  $\forall i \in \mathcal{I}$ .

With the constructed policies  $\{\bar{\sigma}^\ell, \ell = 1, 2, \dots\}$ , we are ready to define the policy  $\bar{\sigma}$ . We let  $\bar{\sigma}$  choose an index  $\ell$  with probability  $p^\ell$  at time 0, and then follow  $\bar{\sigma}^\ell$ . We set the probability  $p^\ell$  as

$$p^\ell = \frac{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})}{\sum_{\ell'=1}^{\infty} \mathbb{P}(K^{\ell'}(\infty) \neq \mathbf{0})} = \frac{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})}{\sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}(\infty)]}, \quad \forall \ell = 1, 2, \dots \quad (31)$$

where the second inequality uses the fact that  $\sum_{\ell'=1}^{\infty} \mathbb{P}(K^{\ell'}(\infty) \neq \mathbf{0}) = \sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}(\infty)]$ . Then under  $\bar{\sigma}$ , we have

$$\begin{aligned} \mathbb{E} [h(\bar{K}(\infty))] &= \sum_{\ell=1}^{\infty} p^\ell \mathbb{E} [h(\bar{K}^\ell(\infty))] \stackrel{(a)}{=} \sum_{\ell=1}^{\infty} \frac{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})}{\sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}(\infty)]} \cdot \frac{\mathbb{E} [h(K^\ell(\infty))]}{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})} \\ &= \frac{\sum_{\ell=1}^{\infty} \mathbb{E} [h(K^\ell(\infty))]}{\sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}(\infty)]} = \frac{\sum_{\mathbf{k} \neq \mathbf{0}} h(\mathbf{k}) \mathbb{E}[X_{\mathbf{k}}(\infty)]}{\sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}(\infty)]} = C(\sigma), \end{aligned}$$

where (a) follows from (30) and (31). Observe that under  $\bar{\sigma}$ ,  $\bar{K}(\infty) \neq \mathbf{0}$  almost surely, we thus have

$$\mathbb{E} [h(\bar{K}(\infty)) | \bar{K}(\infty) \neq \mathbf{0}] = \mathbb{E} [h(\bar{K}(\infty))] = C(\sigma),$$

which proves (27). Moreover, for each  $i \in \mathcal{I}$  the request rate  $\bar{\lambda}_i$  is given by

$$\bar{\lambda}_i = \sum_{\ell=1}^{\infty} p^\ell \cdot \bar{\lambda}_i^\ell = \sum_{\ell=1}^{\infty} \frac{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})}{\sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}(\infty)]} \cdot \frac{\lambda_i^\ell}{\mathbb{P}(K^\ell(\infty) \neq \mathbf{0})} = \frac{\sum_{\ell=1}^{\infty} \lambda_i^\ell}{\sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}(\infty)]} = \frac{\lambda_i r}{N(\sigma)}.$$

This proves (28). By the argument presented at the beginning of the proof, we get  $N(\sigma) \geq \bar{N}^*$ .  $\square$

## B THE REST OF THE PROOFS NEEDED FOR THEOREM 3 (CONVERSION THEOREM)

### B.1 Proof of Lemma 1

PROOF. We will show that under the JRS policy, the Markov chain for the system state (represented as  $((K^\ell(t))_{\ell=1,2,\dots}, \zeta^{1:L}(t), \eta^{1:L}(t)))$  has a unique stationary distribution by first arguing that it is  $\mathbf{k}^0$ -irreducible (here being  $\mathbf{k}^0$ -irreducible means the Markov chain has a state that can be reached

by all other states through transitions, “ $\mathbf{k}^0$ ” in “ $\mathbf{k}^0$ -irreducible” does not refer to any specific states), and then use Foster-Lyapunov theorem to show the positive recurrence [see, e.g., 25]. Combining  $\mathbf{k}^0$ -irreducibility with positive recurrence, we can conclude that the Markov chain under study has a unique stationary distribution.

First, we show that the Markov chain  $((\mathbf{K}^\ell(t))_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  is  $\mathbf{k}^0$ -irreducible. Specifically, observe that the Markov chain starting from any state  $((\mathbf{k}^\ell)_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}, \boldsymbol{\eta}^{1:L})$  can reach the state  $(\mathbf{k}^{1:L}, (\mathbf{0})_{\ell=L+1,\dots}, \mathbf{0}^{1:L}, \mathbf{0}^{1:L})$ , after experiencing a sequence of departures and arrivals that clears up all the tokens, virtual jobs and jobs on backup servers. Further, letting  $\tilde{\mathbf{k}}$  be the configuration reachable by all other configuration in the single-server system under the policy  $\bar{\sigma}$ , we argue that starting from any states of the form  $(\mathbf{k}^{1:L}, (\mathbf{0})_{\ell=L+1,\dots}, \mathbf{0}^{1:L}, \mathbf{0}^{1:L})$ , the Markov chain can reach the state  $(\tilde{\mathbf{k}}^{1:L}, (\mathbf{0})_{\ell=L+1,\dots}, \mathbf{0}^{1:L}, \mathbf{0}^{1:L})$ . Because for any  $\ell \leq L$ , there is a transition path from  $\mathbf{k}^\ell$  to  $\tilde{\mathbf{k}}$ , consider the sequence of events where each  $\mathbf{K}^\ell(t)$  transitions independently following the path, and the jobs arrive right after  $\mathbf{K}^\ell(t)$  making a request, so that the tokens are checked out before  $\mathbf{K}^\ell(t)$  has a further transition. In this way, each  $\mathbf{K}^\ell(t)$  with  $\ell \leq L$  can eventually reach  $\tilde{\mathbf{k}}$  from  $\mathbf{k}^\ell$ . This proves the  $\mathbf{k}^0$ -irreducibility of  $((\mathbf{K}^\ell(t))_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$ .

Next, we show that  $((\mathbf{K}^\ell(t))_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  satisfies the Foster-Lyapunov criterion, i.e.,

$$\widehat{G}g \leq -1 + b\mathbf{1}_{\{S\}}, \quad (32)$$

where  $g$  is a non-negative function of the states,  $S$  is a finite set,  $b$  is a finite number, and  $\widehat{G}$  is the infinitesimal generator of the continuous-time Markov chain. Let  $t_i$  be the expected remaining time in the system when a job is in phase  $i$  for each  $i \in \mathcal{I}$ . According to the job model, we have the recurrence relation

$$\left(\mu_{i\perp} + \sum_{i' \in \mathcal{I} : i' \neq i} \mu_{ii'}\right)t_i = \sum_{i' \in \mathcal{I} : i' \neq i} \mu_{ii'}t_{i'} \quad \forall i \in \mathcal{I}. \quad (33)$$

We construct a Lyapunov function  $g$  as follows:

$$g((\mathbf{k}^\ell)_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}, \boldsymbol{\eta}^{1:L}) = \sum_{i \in \mathcal{I}} \sum_{\ell=1}^{\infty} t_i k_i^\ell + \sum_{i \in \mathcal{I}} \sum_{\ell=1}^{\infty} t_i \zeta_i^\ell. \quad (34)$$

Using the relation (33), it can be verified that the drift of  $g$  satisfies

$$\begin{aligned} & \widehat{G}g((\mathbf{k}^\ell)_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}, \boldsymbol{\eta}^{1:L}) \\ & \leq \sum_{i \in \mathcal{I}} \left( \lambda_i t_i r - \sum_{\ell=1}^{\infty} k_i^\ell \right) + \sum_{i \in \mathcal{I}} \left( \sum_{\ell=1}^L \sum_{(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k}^\ell)} \gamma_{\mathbf{k}^\ell, (\mathbf{k}', \mathbf{a})} a_i t_i - \sum_{\ell=1}^{\infty} \zeta_i^\ell \right) \\ & \leq \sum_{i \in \mathcal{I}} \left( \lambda_i t_i r + L \cdot \max_{\mathbf{k} \in \mathcal{K}} \sum_{(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k})} \gamma_{\mathbf{k}, (\mathbf{k}', \mathbf{a})} a_i t_i \right) - \left( \sum_{\ell=1}^{\infty} \sum_{i \in \mathcal{I}} k_i^\ell + \sum_{\ell=1}^{\infty} \sum_{i \in \mathcal{I}} \zeta_i^\ell \right), \end{aligned}$$

where the first inequality uses the fact that virtual jobs are generated at a rate no faster than the total rate of job requests. Then the Foster-Lyapunov criterion in (32) is satisfied with  $b$  and  $S$  given by

$$\begin{aligned} b &= \sum_{i \in \mathcal{I}} \left( \lambda_i t_i r + L \cdot \max_{\mathbf{k} \in \mathcal{K}} \sum_{(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k})} \gamma_{\mathbf{k}, (\mathbf{k}', \mathbf{a})} a_i t_i \right), \\ S &= \{((\mathbf{k}^\ell)_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}, \boldsymbol{\eta}^{1:L}) : g((\mathbf{k}^\ell)_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}, \boldsymbol{\eta}^{1:L}) \leq b + 1\}. \end{aligned}$$

By the Foster-Lyapunov theorem,  $((\mathbf{K}^\ell(t))_{\ell=1,2,\dots}, \boldsymbol{\zeta}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  is positive recurrent.  $\square$



## B.2 Proofs of Lemma 3 and Lemma 4

In this subsection, we prove Lemma 3 and Lemma 4 together. We begin by introducing some notations. We use  $U^\ell \triangleq (\widehat{K}^\ell, \eta^\ell)$  to represent the state of the  $\ell$ -th server, and use  $U^{1:L}$  to represent the joint state of the first  $L$  servers. We also use the lowercase  $u^\ell, u^{1:L}$  to represent the realizations of the corresponding random variables. We denote the total number of type  $i$  virtual jobs as  $V_i \triangleq \sum_{\ell=1}^L \zeta_i^\ell$  for  $i \in \mathcal{I}$ , and its realization as  $v_i$ . We denote the total number of type  $i$  jobs on backup servers as  $Y_i$  for  $i \in \mathcal{I}$ , and its realizations as  $y_i$ . We also denote the total number of type  $i$  tokens throughout the system as  $Z_i \triangleq \sum_{\ell=1}^L \eta_i^\ell$ , and its realization as  $z_i$ . Our goal can be rewritten as proving  $\mathbb{E}[V_i] = O(\sqrt{r})$  and  $\mathbb{E}[Y_i] = O(\sqrt{r})$  for each  $i \in \mathcal{I}$ .

We first give an overview of the proof. Observe that in our model, the expected time that a job stay in the system is fixed. As a result, bounding the number of virtual jobs or jobs on backup servers in the system is equivalent to bounding the rate that they are generated, according to Little's Law. By our construction of the policy, the rate of generating those jobs is closely related to the dynamics of the total number of type  $i$  tokens  $Z_i(t)$ .

To describe the dynamics of  $Z_i(t)$ , we first introduce two functions  $dv_i$  and  $dy_i$ :

$$\begin{aligned} dv_i(a_i, z_i) &\triangleq (z_i + a_i - \eta_{\max})^+, \\ dy_i(z_i) &\triangleq (1 - z_i)^+. \end{aligned}$$

The function  $dv_i$  represents the increment in the number of type  $i$  virtual jobs due to the event that the total number of type  $i$  tokens on the first  $L$  servers exceeds the token limit  $\eta_{\max}$ . The function  $dy_i$  corresponds to the increment in the total number of type  $i$  jobs on backup servers due to the event that a type  $i$  job arrives to the system without seeing a type  $i$  token. For a function  $g: (\mathcal{K} \times \mathcal{K})^L \rightarrow \mathbb{R}$  that only depends on the number of type  $i$  tokens  $z_i$ , its drift can be written as

$$\begin{aligned} \widehat{G}g(u^{1:L}) &= \sum_{\ell=1}^L \sum_{(k', a) \in E(k^\ell)} \gamma_{k^\ell, (k', a)} (g(z_i + a_i - dv_i(a_i, z_i)) - g(z_i)) \mathbb{1}_{\{\eta^\ell=0\}} \\ &\quad + \lambda_i r (g(z_i - 1 + dy_i(z_i)) - g(z_i)). \end{aligned} \quad (35)$$

We abuse the notation of  $g$  here. For ease of exposition, we will simply write  $dv_i$  and  $dy_i$  to represent  $dv_i(a_i, z_i)$  and  $dy_i(z_i)$ .

By construction, the total number of type  $i$  tokens  $\{Z_i(t)\}$  is a stochastic process constrained within  $[0, \eta_{\max}]$ . Note that  $Z_i(t)$  increases when some servers request new tokens, and decreases when a real or virtual arrival checks out the token or when some servers have the excessive tokens removed. When  $Z_i(t)$  is away from the boundaries, the average rate that it increases is approximately equal to  $\lambda_i r$ , and the average rate that  $Z_i(t)$  decreases is given by

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{(k', a) \in E(\widehat{K}^\ell)} \gamma_{\widehat{K}^\ell, (k', a)} a_i \mathbb{1}_{\{\eta^\ell=0\}} \right] \approx \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{(k', a) \in E(\overline{K}^\ell)} \gamma_{\overline{K}^\ell, (k', a)} a_i \right] = L \cdot \bar{\lambda}_i \approx \lambda_i r,$$

where we have used the approximations that  $\widehat{K}^\ell \stackrel{d}{\approx} \overline{K}^\ell$ ,  $\mathbb{1}_{\{\eta^\ell=0\}} \approx 1$  and  $L = \lceil \bar{N} \rceil \approx \bar{N}$ .

As  $\{Z_i(t)\}$  randomly moves up and down with approximately the same rate and reflects on the boundaries of 0 and  $\eta_{\max}$ , it behaves as a reflected simple symmetric random walk. Intuitively speaking, the steady-state distribution of  $Z_i$  is approximately a uniform distribution over  $[0, \eta_{\max}]$ . Recall that  $dv_i$  and  $dy_i$  can only be non-zero when  $Z_i(t)$  is near the boundaries. Since the length of the interval  $\eta_{\max} = \Theta(\sqrt{r})$ , we can expect that  $dv_i$  and  $dy_i$  diminish as  $r \rightarrow \infty$ .

In the proof, we first establish the relationship between  $\mathbb{E}[V_i]$ ,  $\mathbb{E}[Y_i]$  and  $dv_i$ ,  $dy_i$  using Little's Law. Then we derive bounds on  $dv_i$  and  $dy_i$  by analyzing the drift of several test functions of

$Z_i$ . This step is implicitly based on the intuition that  $Z_i$  is approximately uniformly distributed over  $[0, \eta_{\max}]$ , with the tokens being generated and eliminated at similar rates. Finally, we invoke Lemma 2 to show that the tokens are indeed generated and eliminated at similar speeds, which leads to bounds on  $dv_i$  and  $dy_i$ .

Finally, we make some additional remarks on the notations. First,  $dv_i$  and  $dy_i$  depend on the total number of type  $i$  tokens  $z_i$  and the number of newly requested jobs  $a_i$ , although we omit the dependency expression for ease of exposition. Second, we abuse the notation  $dv_i$  and  $dy_i$  to denote the corresponding random variables. We also write  $\sum_{k',a}$  as a shorthand for  $\sum_{(k',a) \in E(k^\ell)}$  when the context is clear.

**B.2.1 Proofs of Lemma 3 and Lemma 4.** We are now ready to prove Lemma 3 and Lemma 4.

**PROOF. Step 1: Bounding Virtual Jobs and Jobs on Backup Servers using Little's Law.** We first apply Little's Law to  $V_i$  and  $Y_i$ . For each type  $i$ , we let the expected time that a type  $i$  job stays in the system be  $t_i$ . Let  $t_{\max} = \max_{i \in I} t_i$ . Because there are only finitely many types of jobs, and each job spends finite expected time in the system,  $t_{\max}$  is a finite constant. By Little's law,

$$\mathbb{E}[V_i] \leq t_{\max} \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{k',a} Y_{\widehat{K}^\ell, (k',a)} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right], \quad (36)$$

$$\mathbb{E}[Y_i] \leq t_{\max} \mathbb{E} [\lambda_i r \cdot dy_i]. \quad (37)$$

**Step 2: Drift Analysis.** The above two equations (36) and (37) suggest that we can derive upper bounds on  $\mathbb{E}[V_i]$  and  $\mathbb{E}[Y_i]$  by analyzing the following two terms:

- $\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{k',a} Y_{\widehat{K}^\ell, (k',a)} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right]$ , interpreted as the average rate that  $Z_i(t)$  reflects on the boundary at  $\eta_{\max}$ ;
- $\mathbb{E} [\lambda_i r \cdot dy_i]$ , interpreted as the average rate that  $Z_i(t)$  reflects on the boundary at 0.

We establish the relationships of  $dv_i$ ,  $dy_i$  and  $Z_i$  by analyzing the drift of two test functions  $g$ .

Letting  $g(z_i) = z_i$  and taking steady-state expectation over its drift, by (35) and the fact that the drift is zero in steady state, we get

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{k',a} Y_{\widehat{K}^\ell, (k',a)} (a_i - dv_i) \mathbb{1}_{\{\eta^\ell=0\}} + \lambda_i r (-1 + dy_i) \right] = 0. \quad (38)$$

Similarly, letting  $g(z_i) = z_i^2$  and taking steady-state expectation over its drift, one can verify that

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{k',a} Y_{\widehat{K}^\ell, (k',a)} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right] \quad (39)$$

$$= \frac{1}{\eta_{\max}} \cdot \mathbb{E} \left[ \left( \sum_{\ell=1}^L \sum_{k',a} Y_{\widehat{K}^\ell, (k',a)} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \lambda_i r \right) \cdot Z_i \right] \quad (40)$$

$$+ \frac{1}{2\eta_{\max}} \cdot \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{k',a} Y_{\widehat{K}^\ell, (k',a)} (a_i^2 - (dv_i)^2) \mathbb{1}_{\{\eta^\ell=0\}} + \lambda_i r \cdot (1 - (dy_i)^2) \right]. \quad (41)$$

Readers may refer to the complete calculation at the end of this subsection.

**Step 3: Estimating the Terms Obtained from Drift Analysis.** We will first focus on bounding  $\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{k',a} Y_{\widehat{K}^\ell, (k',a)} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right]$  analyzing the two terms in (40) and (41) separately. Then we invoke (38) to bound  $\mathbb{E} [\lambda_i r \cdot dy_i]$ .

The term in (41) is easy to deal with. Observe that the number of jobs requested each time should be no more than the maximal number of jobs that a server can hold, i.e.,  $a_i \leq K_{\max}$ , so

$$\begin{aligned}
 (41) &\leq \frac{1}{2\eta_{\max}} \cdot \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} K_{\max}^2 + \lambda_i r \right] \\
 &\leq \frac{1}{2\eta_{\max}} \cdot \mathbb{E} \left[ \sum_{\ell=1}^L \gamma_{\max} K_{\max}^2 + \lambda_i r \right] \\
 &= O(\sqrt{r}), \tag{42}
 \end{aligned}$$

where in the second inequality we have used the fact that the total rate is uniformly bounded by  $\gamma_{\max}$ , and the last step uses the facts that  $L = O(r)$  and  $\eta_{\max} = \Theta(\sqrt{r})$ .

To bound the term in (40), first observe that  $Z_i \leq \eta_{\max}$ , which implies that

$$(40) \leq \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \lambda_i r \right\|.$$

The term on the RHS of the above equation is the expected absolute difference between the rates of generating and eliminating type  $i$  tokens, which can be shown to be small relative to  $r$ . Specifically, we claim that

$$\mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \lambda_i r \right\| = O(\sqrt{r}). \tag{43}$$

To show (43), first notice that we can remove the indicator  $\mathbb{1}_{\{\eta^\ell=0\}}$  without introducing much error:

$$\begin{aligned}
 &\mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \lambda_i r \right\| \\
 &\leq \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i - \lambda_i r \right\| + \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i \mathbb{1}_{\{\eta^\ell \neq 0\}} \right\| \\
 &\leq \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i - \lambda_i r \right\| + \mathbb{E} \left\| \sum_{\ell=1}^L \gamma_{\max} K_{\max} \mathbb{1}_{\{\eta^\ell \neq 0\}} \right\| \\
 &\leq \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i - \lambda_i r \right\| + \gamma_{\max} K_{\max} |\mathcal{I}| \eta_{\max},
 \end{aligned}$$

where the first inequality is due to triangular inequality, the second inequality is due to the definition of  $\gamma_{\max}$ , and the last inequality is because  $\sum_{\ell=1}^L \mathbb{1}_{\{\eta^\ell \neq 0\}} \leq |\mathcal{I}| \eta_{\max}$ . It remains to bound the term  $\mathbb{E} \left\| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i - \lambda_i r \right\|$ , which can be seen as showing that the rate of generating type  $i$  tokens concentrates around the type  $i$  jobs' arrival rate  $\lambda_i r$ . It is natural to think of using some Law of Large Numbers. Unfortunately,  $\sum_{\ell=1}^L \sum_{k',a} \gamma_{\widehat{K}^\ell, (k',a)} a_i$  is not a sum of i.i.d. random variables due to dependencies among  $\widehat{K}^\ell$  for different  $\ell$ 's. As a result, we want to invoke the Wasserstein distance bound in Lemma 2 to replace  $\widehat{K}^\ell$  in the above expression with  $\overline{K}^\ell$ . We define the function  $f(\mathbf{k}^{1:L})$  as

$$f(\mathbf{k}^{1:L}) = \frac{1}{2\gamma_{\max} K_{\max}} \left| \sum_{\ell=1}^L \sum_{k',a} \gamma_{\mathbf{k}^\ell, (k',a)} a_i - \lambda_i r \right|. \tag{44}$$

We claim that  $f \in \text{Lip}(1)$ . For any two  $\mathbf{k}^{1:L,1}, \mathbf{k}^{1:L,2}$ ,

$$\begin{aligned}
& 2\gamma_{\max}K_{\max} \cdot \left( f(\mathbf{k}^{1:L,1}) - f(\mathbf{k}^{1:L,2}) \right) \\
&= \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\mathbf{k}^{\ell,1},(\mathbf{k}',a)} a_i - \lambda_i r \right| - \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\mathbf{k}^{\ell,2},(\mathbf{k}',a)} a_i - \lambda_i r \right| \\
&\leq \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\mathbf{k}^{\ell,1},(\mathbf{k}',a)} a_i - \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\mathbf{k}^{\ell,2},(\mathbf{k}',a)} a_i \right| \\
&= \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} (\gamma_{\mathbf{k}^{\ell,1},(\mathbf{k}',a)} - \gamma_{\mathbf{k}^{\ell,2},(\mathbf{k}',a)}) a_i \mathbb{1}_{\{(\mathbf{k}^{\ell,1}) \neq (\mathbf{k}^{\ell,2})\}} \right| \\
&\leq \sum_{\ell=1}^L \sum_{\mathbf{k}',a} |\gamma_{\mathbf{k}^{\ell,1},(\mathbf{k}',a)} - \gamma_{\mathbf{k}^{\ell,2},(\mathbf{k}',a)}| \cdot K_{\max} \cdot \|\mathbf{k}^{\ell,1} - \mathbf{k}^{\ell,2}\| \\
&\leq \sum_{\ell=1}^L 2\gamma_{\max}K_{\max} \cdot \|\mathbf{k}^{\ell,1} - \mathbf{k}^{\ell,2}\| \\
&= 2\gamma_{\max}K_{\max} \cdot \|\mathbf{k}^{1:L,1} - \mathbf{k}^{1:L,2}\|,
\end{aligned}$$

where the first inequality is due to triangular inequality; the second inequality uses the fact that  $a_i \leq K_{\max}$  and  $\mathbb{1}_{\{(\mathbf{k}^{1,\ell} \neq \mathbf{k}^{2,\ell})\}} \leq \|\mathbf{k}^{1,\ell} - \mathbf{k}^{2,\ell}\|$ ; the third inequality uses triangular inequality, the fact that the total rate at a configuration  $\mathbf{k}$  is bounded by  $\gamma_{\max}$  and the property of the  $L^1$  norm  $\|\cdot\|$ . Therefore,  $f \in \text{Lip}(1)$ . The Lipschitz continuity of  $f$  allows us to invoke Lemma 2 and get

$$\mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\widehat{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i - \lambda_i r \right| \right] - \mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i - \lambda_i r \right| \right] \leq 2\gamma_{\max}K_{\max} \cdot O(\sqrt{r}).$$

Therefore,

$$\mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\widehat{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i - \lambda_i r \right| \right] \leq \mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i - \lambda_i r \right| \right] + O(\sqrt{r}). \quad (45)$$

Observe that under a Markovian policy, the request rate of type  $i$  jobs can be written as  $\bar{\lambda}_i = \mathbb{E}[\sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i]$ , so we have

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i - \lambda_i r \right] = \bar{\lambda}_i \cdot \lceil \bar{N} \rceil - \lambda_i r = O(1). \quad (46)$$

Moreover, because  $\sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i$  are i.i.d. for  $\ell = 1, \dots, L$ , we have

$$\begin{aligned}
\mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i - \lambda_i r \right| \right] &\leq \sqrt{\mathbb{E} \left[ \left( \sum_{\ell=1}^L \sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i - \lambda_i r \right)^2 \right]} + O(1) \\
&= \sqrt{\sum_{\ell=1}^L \text{Var} \left( \sum_{\mathbf{k}',a} \gamma_{\overline{\mathbf{K}}^{\ell},(\mathbf{k}',a)} a_i \right)} + O(1) \\
&= O(\sqrt{r}).
\end{aligned}$$

Therefore, by combining the arguments above, we get

$$\begin{aligned}
\mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \lambda_i r \right| \right] &\leq \mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i - \lambda_i r \right| \right] + O(\sqrt{r}) \\
&\leq \mathbb{E} \left[ \left| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i - \lambda_i r \right| \right] + O(\sqrt{r}) \\
&\leq O(\sqrt{r}), \tag{47}
\end{aligned}$$

which proves (43). This implies that the term in (40) is also in  $O(\sqrt{r})$ .

Combining the bounds on the terms in (40) and (41), we get

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right] = O(\sqrt{r}).$$

Finally, we bound  $\mathbb{E}[\lambda_i r \cdot dy_i]$ . We rearrange the terms in (38) and get

$$\begin{aligned}
\mathbb{E}[\lambda_i r \cdot dy_i] &= \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} (-a_i + dv_i) \mathbb{1}_{\{\eta^\ell=0\}} + \lambda_i r \right] \\
&= \mathbb{E} \left[ - \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \mathbb{1}_{\{\eta^\ell=0\}} + \lambda_i r \right] + O(\sqrt{r}).
\end{aligned}$$

By (43), we have  $\mathbb{E} \left[ - \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \mathbb{1}_{\{\eta^\ell=0\}} + \lambda_i r \right] = O(\sqrt{r})$ . Therefore,

$$\mathbb{E}[\lambda_i r \cdot dy_i] = O(\sqrt{r}).$$

We invoke the equations (36) and (37) that we get at the beginning of the proof, and conclude that

$$\mathbb{E}[V_i] \leq t_{\max} \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right] = O(\sqrt{r}).$$

$$\mathbb{E}[Y_i] \leq t_{\max} \mathbb{E}[\lambda_i r \cdot dy_i] = O(\sqrt{r}).$$

This finishes the proof.  $\square$

**B.2.2 Deriving the equality in (39).** We show the calculation detail of deriving the following equality.

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right] \tag{39}$$

$$= \frac{1}{\eta_{\max}} \cdot \mathbb{E} \left[ \left( \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \lambda_i r \right) \cdot Z_i \right] \tag{40}$$

$$+ \frac{1}{2\eta_{\max}} \cdot \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} (a_i^2 - (dv_i)^2) \mathbb{1}_{\{\eta^\ell=0\}} + \lambda_i r \cdot (1 - (dy_i)^2) \right]. \tag{41}$$

The equality is obtained by considering the drift of the function  $g(z_i) = z_i^2$ , which is zero in steady state. Recall that the drift of  $g(z_i)$  is given by

$$\begin{aligned} \widehat{G}g(z_i, \mathbf{u}^{1:L}) &= \sum_{\ell=1}^L \sum_{\mathbf{k}', a} \gamma_{\mathbf{k}', (\mathbf{k}', a)} (g(z_i + a_i - dv_i) - g(z_i)) \mathbb{1}_{\{\eta^\ell=0\}} \\ &\quad + \lambda_i r (g(z_i - 1 + dy_i) - g(z_i)). \end{aligned} \quad (35)$$

We will first calculate  $g(z_i + a_i - dv_i) - g(z_i)$ , then  $g(z_i - 1 + dy_i) - g(z_i)$ , and finally plug them into the (35).

The calculation of  $g(z_i + a_i - dv_i) - g(z_i)$  utilizes the following property of  $dv_i$ :

$$(z_i + a_i - dv_i) \cdot dv_i = \eta_{\max} \cdot dv_i. \quad (48)$$

This property follows from the definition  $dv_i = (z_i + a_i - \eta_{\max})^+$ . Intuitively, this is because  $dv_i$  is the “force” that pushes  $z_i$  back when it hits the boundary at  $\eta_{\max}$ . Using the property, we have

$$\begin{aligned} &g(z_i + a_i - dv_i) - g(z_i) \\ &= (z_i + a_i - dv_i)^2 - z_i^2 \\ &= (z_i + a_i - dv_i)^2 - (z_i + a_i - dv_i - a_i + dv_i)^2 \\ &= (z_i + a_i - dv_i)^2 - ((z_i + a_i - dv_i)^2 + 2(-a_i + dv_i) \cdot (z_i + a_i - dv_i) + (-a_i + dv_i)^2) \\ &= 2(a_i - dv_i) \cdot (z_i + a_i - dv_i) - (-a_i + dv_i)^2 \\ &= 2a_i \cdot (z_i + a_i - dv_i) - (-a_i + dv_i)^2 - 2dv_i \cdot \eta_{\max} \\ &= 2a_i \cdot z_i + a_i^2 - (dv_i)^2 - 2dv_i \cdot \eta_{\max}. \end{aligned}$$

The second last equality is due to (48), and the rest are all algebraic manipulations.

We carry out a similar calculation for  $g(z_i - 1 + dy_i) - g(z_i)$ :

$$\begin{aligned} &g(z_i - 1 + dy_i) - g(z_i) \\ &= 2(-1 + dy_i) \cdot z_i + (-1 + dy_i)^2 \\ &= -2z_i + 2z_i \cdot dy_i + 1 - 2dy_i + (dy_i)^2 \\ &= -2z_i + 1 + 2(z_i - 1 + dy_i) \cdot dy_i - (dy_i)^2 \\ &= -2z_i + 1 - (dy_i)^2. \end{aligned}$$

where the last equality is due to the property that

$$(z_i - 1 + dy_i) \cdot dy_i = 0, \quad (49)$$

and the rest are all algebraic manipulations.

Putting together,

$$\begin{aligned} \widehat{G}g(z_i, \mathbf{u}^{1:L}) &= \sum_{\ell=1}^L \sum_{\mathbf{k}', a} \gamma_{\mathbf{k}', (\mathbf{k}', a)} (g(z_i + a_i - dv_i) - g(z_i)) \mathbb{1}_{\{\eta^\ell=0\}} \\ &\quad + \lambda_i r (g(z_i - 1 + dy_i) - g(z_i)) \\ &= \sum_{\ell=1}^L \sum_{\mathbf{k}', a} \gamma_{\mathbf{k}', (\mathbf{k}', a)} (2a_i \cdot z_i + a_i^2 - (dv_i)^2 - 2dv_i \cdot \eta_{\max}) \mathbb{1}_{\{\eta^\ell=0\}} \\ &\quad + \lambda_i r \cdot (-2z_i + 1 - (dy_i)^2). \end{aligned}$$



After recombining the terms, we get

$$\begin{aligned}
& \eta_{\max} \cdot \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right] \\
&= \mathbb{E} \left[ \left( \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \lambda_i r \right) \cdot Z_i \right] \\
&+ \frac{1}{2} \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \gamma_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} (a_i^2 - (dv_i)^2) \mathbb{1}_{\{\eta^\ell=0\}} + \lambda_i r \cdot (1 - (dy_i)^2) \right]
\end{aligned}$$

This finishes the calculation.

### B.3 Proof of Conversion Theorem without Assuming Irreducibility

In this subsection, we prove Theorem 3 without assuming  $\mathbf{k}^0$ -irreducibility of the subroutine  $\bar{\sigma}$ . Specifically, suppose that we have a Markovian single-server policy  $\bar{\sigma}$  and an initial distribution  $p_j$  over its recurrent classes  $S_j$  for  $j = 1, 2, \dots, J$ , we will construct an infinite-server policy  $\sigma$  such that (5) and (6) still hold. The basic idea is to decompose a general Markovian single-server policy  $\bar{\sigma}$  into multiple  $\mathbf{k}^0$ -irreducible Markovian policies, each induces one recurrent class and preserves stationary distribution and the throughput of  $\bar{\sigma}$  on that recurrent class, as stated in Lemma 5 below.

**Lemma 5** (Decomposing The Reducible Policy). *Let  $\bar{\sigma}$  be a general single-server Markovian policy with recurrent classes  $S_j$  for  $j = 1, 2, \dots, J$ . Then for each  $j$  exists a Markovian policy  $\bar{\sigma}^j$  such that*

- *The induced Markov chain is  $\mathbf{k}^0$ -irreducible with the unique recurrent class being  $S_j$ ;*
- *The stationary distribution is the same as the stationary distribution under  $\bar{\sigma}$  starting from a configuration in  $S_j$ .*

**PROOF.** For each  $j = 1, 2, \dots, J$ , we define the policy  $\bar{\sigma}^j$  as follows: when the system has configuration  $\mathbf{k} \in S_j$ , the policy  $\bar{\sigma}^j$  makes the same decisions as  $\bar{\sigma}$ ; when  $\mathbf{k} \notin S_j$  and  $\mathbf{k} = \mathbf{0}$ , the policy starts a timer whose duration follows an exponential distribution with rate 1 and immediately adds  $\mathbf{k}^0$  many jobs of each type when the timer ticks for some arbitrary  $\mathbf{k}^0 \in S_j$ ; when  $\mathbf{k} \notin S_j$  and  $\mathbf{k} \neq \mathbf{0}$ , the policy does not request any jobs.

We show that under the new policy  $\bar{\sigma}^j$ ,  $S_j$  is also a recurrent class of the induced Markov chain. This is because if the system starts from a configuration  $\mathbf{k} \in S_j$ , then it will stay in  $S_j$  since it makes the same decisions and has the same transitions as under the policy  $\bar{\sigma}$ . Because  $S_j$  is a recurrent class under  $\bar{\sigma}$ , it is still a recurrent class under  $\bar{\sigma}^j$ .

To show that the Markov chain induced by  $\bar{\sigma}^j$  is  $\mathbf{k}^0$ -irreducible, observe that starting from any  $\mathbf{k} \notin S_j$ , the system state will return to  $S_j$ . Specifically,

- If the system starts from a configuration  $\mathbf{k}$  such that  $\mathbf{k} \notin S_j$  and  $\mathbf{k} \neq \mathbf{0}$ , then no new jobs will be requested until either  $\mathbf{k} \in S_j$  or  $\mathbf{k} = \mathbf{0}$ . In the latter case, by the construction of the policy, the system jumps to a configuration in  $S_j$  after the next transition.
- If the system starts from a configuration  $\mathbf{k}$  such that  $\mathbf{k} \notin S_j$  and  $\mathbf{k} = \mathbf{0}$ , the system jumps to a configuration in  $S_j$  after the next transition.

The claim that the stationary distribution under  $\bar{\sigma}^j$  is the same as the stationary distribution under  $\bar{\sigma}$  starting from a configuration in  $S_j$  is trivial to show, because when the system initializes from any configuration in  $S_j$ , it stays in  $S_j$  and the transitions are exactly the same under the two policies.  $\square$

The JRS policy with a general Markovian single-server policy as its subroutine  $\bar{\sigma}$  is constructed using the  $\mathbf{k}^0$ -irreducible policies  $\bar{\sigma}^j$ 's obtained from the decomposition of  $\bar{\sigma}$  and the probabilities  $p^j$ 's.

- (1) We divide all the servers into  $J$  server pools, each with infinitely many servers. Let the  $j$ -th server pool run the JRS policy with subroutine  $\bar{\sigma}^j$  (defined in Section 4.2) for each  $j = 1, \dots, J$ .
- (2) Whenever we see an arrival of type  $i$ , we route the job to the  $j$ -th infinite-server system with probability  $\frac{p^j \bar{\lambda}_i^j}{\sum_j p^j \bar{\lambda}_i^j}$  for each  $j = 1, \dots, J$ .

To analyze the policy  $\sigma$ , let  $\pi^j$  and  $\bar{\lambda}_i^j$ 's be the stationary distribution and throughput of the policy  $\bar{\sigma}^j$  for  $j = 1, \dots, J$ . By Lemma 5, we have the following relationships:

$$\pi(\mathbf{k}) = \sum_{j=1}^J p^j \pi^j(\mathbf{k}), \quad \forall \mathbf{k} \in \mathcal{K}, \quad (50)$$

$$\bar{\lambda}_i = \sum_{j=1}^J p^j \bar{\lambda}_i^j, \quad \forall i \in \mathcal{I}. \quad (51)$$

Based on the above relationships, we can prove the general version of the Conversion Theorem that does not require  $\mathbf{k}^0$ -irreducibility.

**PROOF OF THEOREM 3.** For each  $j = 1, \dots, J$ , Lemma 5 implies that the policy  $\bar{\sigma}^j$  and stationary distribution  $\pi^j$  form a feasible solution to the single-server problem  $\bar{\mathcal{P}}((p^j \bar{\lambda}_i^j \bar{N})_{i \in \mathcal{I}}, \epsilon)$ , and the corresponding objective value is  $p^j \bar{N}$ . Consider the infinite-server system with arrival rates  $(p^j \bar{\lambda}_i^j \bar{N})_{i \in \mathcal{I}}$  and budget  $\epsilon$ . As we have proved Theorem 3 for the JRS policy with  $\mathbf{k}^0$ -irreducible subroutine  $\bar{\sigma}^j$ , it follows that

$$\left| \sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}^j] - \left[ p^j \bar{N} \right] \cdot (1 - \pi^j(\mathbf{0})) \right| = O(\sqrt{r}), \quad (52)$$

$$\left| \sum_{\mathbf{k} \neq \mathbf{0}} h(\mathbf{k}) \mathbb{E}[X_{\mathbf{k}}^j] - \left[ p^j \bar{N} \right] \cdot \sum_{\mathbf{k} \neq \mathbf{0}} h(\mathbf{k}) \pi^j(\mathbf{k}) \right| = O(\sqrt{r}), \quad (53)$$

where  $X_{\mathbf{k}}^j$  is the random variable representing the steady-state number of servers in configuration  $\mathbf{k}$  in the infinite-server system under the JRS policy with subroutine  $\bar{\sigma}^j$ .

By the construction of  $\sigma$ , the arrival rate to the  $j$ -th server pool is equal to  $\frac{p^j \bar{\lambda}_i^j}{\sum_j p^j \bar{\lambda}_i^j} \cdot \lambda_i r = \frac{p^j \bar{\lambda}_i^j}{\bar{\lambda}_i} \cdot \lambda_i r = p^j \bar{\lambda}_i^j \bar{N}$ , where the first equality is due to (51), and the second equality is due to the condition that  $\lambda_i r = \bar{\lambda}_i^j \cdot L$ . Therefore, (52) and (53) hold, and we have

$$\begin{aligned} \left| \sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}] - \bar{N} \cdot (1 - \pi(\mathbf{0})) \right| &= \left| \sum_{j=1}^J \sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}^j] - \sum_{j=1}^J p^j \bar{N} \cdot (1 - \pi^j(\mathbf{0})) \right| \\ &\leq \sum_{j=1}^J \left| \sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}^j] - \left[ p^j \bar{N} \right] \cdot (1 - \pi^j(\mathbf{0})) \right| + O(1) \\ &= O(\sqrt{r}). \end{aligned}$$

Here we use (52) and the relationship between  $\pi(\mathbf{k})$  and  $\pi^j(\mathbf{k})$ . Similarly,

$$\left| \sum_{\mathbf{k} \neq \mathbf{0}} h(\mathbf{k}) \mathbb{E}[X_{\mathbf{k}}] - \bar{N} \cdot \sum_{j=1}^J \sum_{\mathbf{k} \neq \mathbf{0}} h(\mathbf{k}) \pi(\mathbf{k}) \right| = \left| \sum_{j=1}^J \sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}^j] - \sum_{j=1}^J p^j \bar{N} \cdot \sum_{\mathbf{k} \neq \mathbf{0}} h(\mathbf{k}) \pi^j(\mathbf{k}) \right|$$

$$\begin{aligned}
&= \sum_{j=1}^J \left| \sum_{\mathbf{k} \neq \mathbf{0}} \mathbb{E}[X_{\mathbf{k}}^j] - \lceil p^j \bar{N} \rceil \cdot \sum_{\mathbf{k} \neq \mathbf{0}} h(\mathbf{k}) \pi^j(\mathbf{k}) \right| + O(1) \\
&= O(\sqrt{r}).
\end{aligned}$$

After re-indexing the servers, we get (5) and (6). The bounds on  $N(\sigma)$  and  $C(\sigma)$ , (7) and (8), follow from (5) and (6). They can be verified in the same way as that in the proof for the irreducible case, so we omit the argument here.  $\square$

## C SOLVING THE SINGLE-SERVER PROBLEM

In this section, we show the equivalence of the single-server problem in (3) and a linear program (LP) as stated in Theorem 4. The equivalence needs to be proved in two directions. In Appendix C.1, we first derive the linear program (61) as a relaxation of the single-server problem (3) so that the optimal value of the LP is a lower bound to the optimal value of (3). Then in Appendix C.2, we will construct a Markovian single-server policy that achieves the optimal value of the LP, which implies the optimality of the policy.

### C.1 Lower Bound via LP Relaxation

In this subsection we derive an LP relaxation of the optimization problem (3), restated below:

$$\begin{aligned}
&\underset{\bar{N}, \bar{\sigma}, \pi}{\text{minimize}} && \bar{N} \\
&\text{subject to} && \mathbb{E} [h(\bar{K}(\infty)) | \bar{K}(\infty) \neq \mathbf{0}] \leq \epsilon, \\
&&& \bar{N} \cdot \bar{\lambda}_i = \lambda_i r, \quad \forall i \in \mathcal{I}.
\end{aligned} \tag{3 revisited}$$

Here we allow  $\bar{\sigma}$  to be non-Markovian, i.e., it can make decisions based on the history, but we still require its performance metrics used in the objective and the constraints of (3) to have well-defined steady-state distributions.

Observe that both  $\bar{K}(\infty)$  and  $\bar{\lambda}_i$  depend on the stationary distribution  $\pi$ , but the constraints in terms of  $\pi$  are implicit. To derive an LP relaxation, we give an explicit characterization of the constraints that must be satisfied by the stationary distribution  $\pi$  induced by any feasible policy  $\bar{\sigma}$ .

To do this, we derive a version of the stationary equation in terms of a quantity called *transition frequency*. The transition frequency of type  $i$  jobs is a function  $u_i: \mathcal{K} \rightarrow \mathbb{R}$  describing the steady-state frequency of requesting a type  $i$  job when the system has configuration  $\mathbf{k}$ . To rigorously define transition frequency, we first introduce a concept called *nominal transition*.

**Definition 1** (Nominal Transition). Consider a single-server system under any policy. When the configuration  $\bar{K}(t)$  transitions from  $\mathbf{k}$  to  $\mathbf{k}' + \mathbf{a}$  for some  $\mathbf{k}, \mathbf{k}' \in \mathcal{K}$  with  $\mathbf{a} = (a_i)_{i \in \mathcal{I}}$  new jobs added into service, we decompose the transition by adding intermediate configurations as illustrated below, where  $\mathbf{k}$  first goes to  $\mathbf{k}'$  if  $\mathbf{k}' \neq \mathbf{k}$ , then add jobs of each type one by one.

$$\mathbf{k} \rightarrow \mathbf{k}' \rightarrow (\mathbf{k}' + \mathbf{e}_{i_1}) \rightarrow \cdots \rightarrow (\mathbf{k}' + a_{i_1} \mathbf{e}_{i_1}) \rightarrow \cdots \rightarrow (\mathbf{k}' + a_{i_1} \mathbf{e}_{i_1} + \cdots + a_{i_{|\mathcal{I}|}} \mathbf{e}_{i_{|\mathcal{I}|}}),$$

where  $(i_1, i_2, \dots, i_{|\mathcal{I}|})$  is a fixed ordering of the set of phases  $\mathcal{I}$ . We call each short transition in the diagram a *nominal transition*.

For  $\mathbf{k}^1, \mathbf{k}^2 \in \mathcal{K}$ , we denote  $F(\mathbf{k}^1, \mathbf{k}^2, t)$  as the cumulative number of nominal transitions from  $\mathbf{k}^1$  to  $\mathbf{k}^2$  during the time interval  $[0, t]$ , which is a random variable with a distribution depending on the single-server policy and initial distribution of configurations.

Note that for any  $i \in \mathcal{I}$  and  $\mathbf{k} \in \mathcal{K}$  s.t.  $k_i \geq 1$ ,  $F(\mathbf{k}, \mathbf{k} - \mathbf{e}_i, t)$  counts the number of times that a type  $i$  job departs when being in configuration  $\mathbf{k}$ . As a result,

$$F(\mathbf{k}, \mathbf{k} - \mathbf{e}_i, t) = \mathcal{N} \left( \int_0^t k_i \mu_{i\perp} \mathbb{1}_{\{K(s)=\mathbf{k}\}} ds \right),$$

where  $\mathcal{N}(t)$  denotes a unit rate Poisson process. If we take expectation, divide both sides by  $t$ , and let  $t \rightarrow \infty$ , we have

$$\lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E}[F(\mathbf{k}, \mathbf{k} - \mathbf{e}_i, t)] = k_i \mu_{i\perp} \cdot \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t \mathbb{P}(K(s) = \mathbf{k}) ds = k_i \mu_{i\perp} \pi(\mathbf{k}). \quad (54)$$

Similarly,  $F(\mathbf{k}, \mathbf{k} - \mathbf{e}_i + \mathbf{e}_{i'}, t)$  counts the number of times a job in phase  $i$  transitions to phase  $i'$  when being in configuration  $\mathbf{k}$  for any  $i, i' \in \mathcal{I}$ ,  $\mathbf{k} \in \mathcal{K}$  s.t.  $i' \neq i$  and  $k_i \geq 1$ , so

$$\lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E}[F(\mathbf{k}, \mathbf{k} - \mathbf{e}_i + \mathbf{e}_{i'}, t)] = k_i \mu_{ii'} \pi(\mathbf{k}). \quad (55)$$

We define transition frequency as follows.

**Definition 2** (Transition Frequency). Transition frequency of type  $i$  jobs at state  $\mathbf{k}$  is the long-run average number of nominal transitions from configuration  $\mathbf{k}$  to  $\mathbf{k} + \mathbf{e}_i$  per unit time,

$$u_i(\mathbf{k}) \triangleq \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)]. \quad (56)$$

The transition frequencies allow us to derive the following version of the stationary equation.

**Lemma 6** (Stationary Equation). *Under any policy, the stationary distribution  $\pi$  and the transition frequency  $u_i$  satisfy the following equation:*

$$\begin{aligned} & \sum_i u_i(\mathbf{k} - \mathbf{e}_i) \mathbb{1}_{\{k_i \geq 1\}} + \sum_i (k_i + 1) \mu_{i\perp} \pi(\mathbf{k} + \mathbf{e}_i) \\ & + \sum_{i, i': i \neq i'} (k_i + 1) \mu_{ii'} \pi(\mathbf{k} + \mathbf{e}_i - \mathbf{e}_{i'}) \mathbb{1}_{\{k_{i'} \geq 1\}} \\ & = \sum_i u_i(\mathbf{k}) + \left( \sum_i k_i \mu_{i\perp} + \sum_{i, i': i \neq i'} k_i \mu_{ii'} \right) \pi(\mathbf{k}) \end{aligned} \quad (57)$$

for any state  $\mathbf{k} \in \mathcal{K}$ , and  $\sum_i, \sum_{i, i'}$  are shorthand for  $\sum_{i \in \mathcal{I}}, \sum_{i, i' \in \mathcal{I}}$ .

**PROOF.** For each configuration  $\mathbf{k} \in \mathcal{K}$ , if we look at the difference between the number of nominal transitions into configuration  $\mathbf{k}$  and that out of configuration  $\mathbf{k}$  by time  $t$ , we have the following equation,

$$\begin{aligned} & \mathbb{1}_{\{\bar{K}(t)=\mathbf{k}\}} - \mathbb{1}_{\{\bar{K}(0)=\mathbf{k}\}} \\ & = \sum_i F(\mathbf{k} - \mathbf{e}_i, \mathbf{k}, t) \mathbb{1}_{\{k_i \geq 1\}} + \sum_i F(\mathbf{k} + \mathbf{e}_i, \mathbf{k}, t) + \sum_{i, i': i \neq i'} F(\mathbf{k} + \mathbf{e}_i - \mathbf{e}_{i'}, \mathbf{k}, t) \mathbb{1}_{\{k_{i'} \geq 1\}} \\ & - \sum_i F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t) - \sum_i F(\mathbf{k}, \mathbf{k} - \mathbf{e}_i, t) \mathbb{1}_{\{k_i \geq 1\}} - \sum_{i, i': i \neq i'} F(\mathbf{k}, \mathbf{k} - \mathbf{e}_i + \mathbf{e}_{i'}, t) \mathbb{1}_{\{k_i \geq 1\}} \end{aligned} \quad (58)$$

By Definition 2 and (54)-(55), if we divide both sides of (58) by  $t$  and let  $t \rightarrow \infty$ , we get the stationary equation in (57).  $\square$

Since the stationary equation in (57) is linear in  $u_i(\mathbf{k})$  and  $\pi(\mathbf{k})$ , we can write it in matrix form:

$$A\pi + \sum_{i \in \mathcal{I}} B_i \mathbf{u}_i = 0, \quad (59)$$

where  $\pi$  and  $u_i$  are column vectors representing  $\pi(\cdot)$ ,  $u_i(\cdot)$ , and  $A, B_i$  are matrices that make (59) equivalent to (57). Therefore, the following three conditions are necessary for any tuple  $(\pi, (u_i)_{i \in \mathcal{I}})$  to be a possible pair of stationary distribution and transition frequencies for a Markovian policy.

$$\begin{aligned} A\pi + \sum_{i \in \mathcal{I}} B_i u_i &= 0 \\ \sum_{\mathbf{k}} \pi(\mathbf{k}) &= 1 \\ \pi, u_i &\geq 0, \quad \forall i \in \mathcal{I} \end{aligned} \quad (60)$$

Based on the characterization of stationary  $\pi$  and  $u_i$ 's in (60), we can now formulate a linear program  $\overline{\mathcal{LP}}((\lambda_i)_{i \in \mathcal{I}}, \epsilon)$ . The linear program has decision variables  $\Phi \in \mathbb{R}$ ,  $\pi \in \mathbb{R}^{\mathcal{K}}$ , and  $u_i \in \mathbb{R}^{\mathcal{K}}$  for  $i \in \mathcal{I}$ , where  $\Phi$  is a factor that scales the throughput of each type of jobs in the direction of  $(\lambda_i)_{i \in \mathcal{I}}$ .

$$\begin{aligned} &\text{maximize} \quad \Phi \\ &\Phi, \pi, (u_i)_{i \in \mathcal{I}} \\ &\text{subject to} \quad \mathbf{h}^T \pi \leq \epsilon(1 - \pi(\mathbf{0})) \\ &\quad \mathbf{1}_0^T u_i = \Phi \cdot \lambda_i \quad \forall i \\ &\quad A\pi + \sum_{i \in \mathcal{I}} B_i u_i = \mathbf{0} \\ &\quad \mathbf{1}^T \pi = 1 \\ &\quad \pi \geq 0, u_i \geq 0 \quad \forall i \in \mathcal{I} \end{aligned} \quad (61)$$

where  $\mathbf{h}$  is the vector form of the cost rate function  $h$ ;  $\mathbf{1}_0^T$  is a  $|\mathcal{K}|$ -dimensional vector with one in all entries except those with  $\sum_{i \in \mathcal{I}} k_i = K_{\max}$ . In addition to the last three constraints on stationarity, the first constraint of  $\overline{\mathcal{LP}}((\lambda_i)_{i \in \mathcal{I}}, \epsilon)$  is the resource contention constraint; the second constraint is because the transition frequency from  $\mathbf{k}$  to  $\mathbf{k} + \mathbf{e}_i$  is equal to the throughput of type  $i$  jobs.

By the construction of  $\overline{\mathcal{LP}}((\lambda_i)_{i \in \mathcal{I}}, \epsilon)$ , any feasible solution of  $\overline{\mathcal{P}}((\lambda_i)_{i \in \mathcal{I}}, \epsilon)$  can be converted to a feasible solution of  $\overline{\mathcal{LP}}((\lambda_i)_{i \in \mathcal{I}}, \epsilon)$ . Let  $\overline{N}^*$  be the optimal value of (3) and  $\Phi^*$  be the optimal value of (61). Then we have

$$\overline{N}^* \geq \frac{r}{\Phi^*}. \quad (62)$$

## C.2 Policy Construction

In this subsection, we describe a procedure that allows us to construct a policy that achieves the lower bound given by the LP relaxation in (61). Specifically, given a feasible solution  $(\pi, (u_i)_{i \in \mathcal{I}})$  to (61), we define an *LP-based policy* that requests jobs as follows:

- Case 1: When the system enters a configuration  $\mathbf{k}$  with  $\pi(\mathbf{k}) \neq 0$ , for each  $i \in \mathcal{I}$ , the policy starts a timer whose duration follows an exponential distribution with rate  $\frac{u_i(\mathbf{k})}{\pi(\mathbf{k})}$ . The policy requests a type  $i$  job when the  $i$ -th timer ticks. When the configuration changes, all timers are canceled.
- Case 2: When the system enters a configuration  $\mathbf{k}$  with  $\pi(\mathbf{k}) = 0$  and  $\sum_{i'} u_{i'}(\mathbf{k}) \neq 0$ , the policy immediately requests a type  $i$  job with probability  $\frac{u_i(\mathbf{k})}{\sum_{i'} u_{i'}(\mathbf{k})}$ .
- Case 3: When the system enters a configuration  $\mathbf{k}$  with  $\pi(\mathbf{k}) = 0$  and  $\sum_{i'} u_{i'}(\mathbf{k}) = 0$ , the policy does not request any jobs.

We denote the LP-based policy based on the solution  $(\pi, (u_i)_{i \in \mathcal{I}})$  as  $\overline{\sigma}(\pi, (u_i)_{i \in \mathcal{I}})$ .

*Remark 2.* Note that the definition of the LP-based policy here is stated from a view different from the view in Section 4.2: here each request only adds one job to the server, and one request can happen immediately after another; while in Section 4.2 each request can add multiple jobs to the server, and there is only one request happening at a time. We refer to the view here as the *impulsive*

view, because multiple requests happening at the same time can be thought of as having an infinite request rate. In contrast, we call the view in Section 4.2 the *non-impulsive view*.

The LP-based policy can be alternatively described using the non-impulsive view if we see multiple requests happening at the same time as one request that adds multiple jobs to the server. More specifically, each reactive request of the LP-based policy is initiated by an internal transition or departure and consists of one or multiple requests of Case 2; each proactive request of the LP-based policy consists of one request in Case 1 and possibly several requests in Case 2.

The following lemma characterizes the steady-state behavior of a single-server system under the LP-based policy.

**Lemma 7** (Properties of LP-based Policies). *Consider a single-server system under the LP-based policy  $\bar{\sigma}(\pi, (\mathbf{u}_i)_{i \in I})$ , where  $(\pi, (\mathbf{u}_i)_{i \in I})$  is a feasible solution to (61). We have that  $\pi$  is a stationary distribution under policy  $\bar{\sigma}$ , and  $(\mathbf{u}_i)_{i \in I}$  are the transition frequencies corresponding to  $\bar{\sigma}$ .*

The proof of Lemma 7 is based on (58), following the same argument as the proof of Lemma 6, as well as an induction argument.

PROOF. Let  $(\pi, (\mathbf{u}_i)_{i \in I})$  be a feasible solution to the LP in (61). To show that  $\pi$  and  $(\mathbf{u}_i)_{i \in I}$  are also the actual stationary distribution and transition frequencies, it suffices to show that if the initial distribution follows  $\pi$ , i.e.

$$\mathbb{P}(\bar{\mathbf{K}}(0) = \mathbf{k}) = \pi(\mathbf{k})$$

then under the policy  $\bar{\sigma}(\pi, (\mathbf{u}_i)_{i \in I})$ , we have

$$\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) dt = \pi(\mathbf{k}), \quad (63)$$

$$\lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, T)] = u_i(\mathbf{k}), \quad (64)$$

where  $F$  is the cumulative number of nominal transitions under the policy  $\bar{\sigma}(\pi, (\mathbf{u}_i)_{i \in I})$  and the initial distribution  $\pi$ .

Our proof is based on the following equation:

$$\begin{aligned} & \frac{d}{dt} \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \Big|_{t=0} \\ &= \sum_i \frac{d}{dt} \mathbb{E}[F(\mathbf{k} - \mathbf{e}_i, \mathbf{k}, t)] \Big|_{t=0} - \sum_i \frac{d}{dt} \mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] \Big|_{t=0} \\ & \quad + \sum_i (k_i + 1) \mu_{i\perp} \pi(\mathbf{k} + \mathbf{e}_i) + \sum_{i,i'} (k_i + 1) \mu_{ii'} \pi(\mathbf{k} + \mathbf{e}_i - \mathbf{e}_{i'}) \mathbb{1}_{\{k_{i'} \geq 1\}} \\ & \quad - \sum_i k_i \mu_{i\perp} \pi(\mathbf{k}) - \sum_{i,i'} k_i \mu_{ii'} \pi(\mathbf{k}), \end{aligned} \quad (65)$$

The equation is a straightforward consequence of (58), following the same argument as the proof of Lemma 6.

We prove the following two equations by induction on  $\sum_{i \in I} k_i$ .

$$\lim_{t \rightarrow 0} \frac{1}{t} \mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] = u_i(\mathbf{k}), \quad (66)$$

$$\frac{d}{dt} \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \Big|_{t=0} = 0. \quad (67)$$



We first consider the base case when  $\sum_{i \in I} k_i = 0$ . In this case,  $\mathbf{k} = \mathbf{0}$  and we have

$$\frac{d}{dt} \mathbb{E} [F(\mathbf{k} - \mathbf{e}_i, \mathbf{k}, t)] \Big|_{t=0} = u_i(\mathbf{k} - \mathbf{e}_i) \mathbb{1}_{\{k_i \geq 1\}} = 0, \quad (68)$$

for all  $i$ . This reduces (65) to

$$\begin{aligned} & \frac{d}{dt} \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \Big|_{t=0} \\ &= \sum_i u(\mathbf{k} - \mathbf{e}_i) \mathbb{1}_{\{k_i \geq 1\}} - \sum_i \frac{d}{dt} \mathbb{E} [F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] \Big|_{t=0} \\ & \quad + \sum_i (k_i + 1) \mu_{i\perp} \pi(\mathbf{k} + \mathbf{e}_i) + \sum_{i,i'} (k_i + 1) \mu_{i,i'} \pi(\mathbf{k} + \mathbf{e}_i - \mathbf{e}_{i'}) \mathbb{1}_{\{k_{i'} \geq 1\}} \\ & \quad - \sum_i k_i \mu_{i\perp} \pi(\mathbf{k}) - \sum_{i,i'} k_i \mu_{i,i'} \pi(\mathbf{k}), \end{aligned} \quad (69)$$

Now we discuss based on whether  $\pi(\mathbf{k}) = 0$ . If  $\pi(\mathbf{k}) \neq 0$ , by the definition of our policy, for all  $i$ ,

$$\frac{d}{dt} \mathbb{E} [F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] \Big|_{t=0} = \frac{u_i(\mathbf{k})}{\pi(\mathbf{k})} \cdot \mathbb{P}(\bar{\mathbf{K}}(0) = \mathbf{k}) = u_i(\mathbf{k}), \quad (70)$$

which is (66). Combining the above equation and the stationary equation (57) satisfied by  $(\pi, (u_i)_{i \in I})$ , we conclude that the RHS of (69) is zero, i.e.

$$\frac{d}{dt} \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \Big|_{t=0} = 0,$$

which is (67). For the case when  $\pi(\mathbf{k}) = 0$  and  $\sum_i u_i(\mathbf{k}) \neq 0$ , because the system immediately leave the configuration  $\mathbf{k}$  after reaching it through a nominal transition,

$$\frac{d}{dt} \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \Big|_{t=0} = 0,$$

i.e., the LHS of (69) is 0. Again we compare (69) against the stationary equation (57) and get

$$\sum_i \frac{d}{dt} \mathbb{E} [F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] \Big|_{t=0} = \sum_i u_i(\mathbf{k}).$$

By the definition of our policy, we have

$$\frac{d}{dt} \mathbb{E} [F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] \Big|_{t=0} = \frac{u_i(\mathbf{k})}{\sum_{i'} u_{i'}(\mathbf{k})} \cdot \sum_{i'} \frac{d}{dt} \mathbb{E} [F(\mathbf{k}, \mathbf{k} + \mathbf{e}_{i'}, t)] \Big|_{t=0} = u_i(\mathbf{k}). \quad (71)$$

which is (66). For the case when  $\pi(\mathbf{k}) = 0$  and  $\sum_i u_i(\mathbf{k}) = 0$ , (57) implies that  $u_i(\mathbf{k} - \mathbf{e}_i) = 0$ ,  $\pi(\mathbf{k} + \mathbf{e}_i) = 0$ ,  $\pi(\mathbf{k} + \mathbf{e}_i - \mathbf{e}_{i'}) = 0$  for any  $i$ . Then (69) is further reduced to

$$\frac{d}{dt} \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \Big|_{t=0} = - \sum_i \frac{d}{dt} \mathbb{E} [F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] \Big|_{t=0}.$$

Because  $\mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \geq 0$  and  $\mathbb{P}(\bar{\mathbf{K}}(0) = \mathbf{k}) = 0$ , the LHS of the above expression is non-negative. However, the RHS of the above expression is non-positive. Therefore, both sides are equal to zero, thus we have  $\frac{d}{dt} \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) \Big|_{t=0} = 0$  and  $\frac{d}{dt} \mathbb{E} [F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] \Big|_{t=0} = u_i(\mathbf{k})$ .

Having proved the base case, we do the induction step. Suppose we have proved (66) and (67) for all  $\mathbf{k}$  such that  $\sum_{i \in I} k_i \leq m - 1$  for some integer  $m \geq 1$ . We consider  $\mathbf{k}$  with  $\sum_{i \in I} k_i = m$ . By the induction hypothesis,

$$\frac{d}{dt} \mathbb{E} [F(\mathbf{k} - \mathbf{e}_i, \mathbf{k}, t)] \Big|_{t=0} = u_i(\mathbf{k} - \mathbf{e}_i) \mathbb{1}_{\{k_i \geq 1\}}. \quad (72)$$

Then we repeat the arguments after (68) of the base case verbatim. By induction, we have proved (66) and (67).

Therefore, given the policy and initial distribution, the distribution of  $\bar{\mathbf{K}}(t)$  is stationary, i.e., we always have  $\mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) = \pi(\mathbf{k})$  for all  $\mathbf{k} \in \mathcal{K}$ . As a result, an analogue of (66) holds for all  $t \geq 0$ :  $\mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)]$  is differentiable with respect to  $t$  and

$$\frac{d}{dt} \mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] = u_i(\mathbf{k}),$$

for all  $\mathbf{k} \in \mathcal{K}$  and all  $i \in \mathcal{I}$ . Therefore,

$$\begin{aligned} \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \mathbb{P}(\bar{\mathbf{K}}(t) = \mathbf{k}) dt &= \lim_{T \rightarrow \infty} \frac{1}{T} \cdot \pi(\mathbf{k}) \cdot T = \pi(\mathbf{k}), \\ \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, T)] &= \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \frac{d}{dt} \mathbb{E}[F(\mathbf{k}, \mathbf{k} + \mathbf{e}_i, t)] dt \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \cdot u_i(\mathbf{k}) \cdot T = u_i(\mathbf{k}). \end{aligned}$$

This completes the proof.  $\square$

### C.3 Proof of Theorem 4

**Theorem 4** (Optimality of Single-OPT). *Given an optimal solution  $(\Phi^*, \pi^*, (\mathbf{u}_i)_{i \in \mathcal{I}})$  to the linear program  $\overline{\mathcal{LP}}((\lambda_i)_{i \in \mathcal{I}}, \epsilon)$ , we can solve the single-server problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$  in (3) to achieve an optimal value  $r/\Phi^*$ , using the optimal policy  $\bar{\sigma}(\pi^*, (\mathbf{u}_i^*)_{i \in \mathcal{I}})$  and the optimal stationary distribution  $\pi^*$ . Moreover, the optimal policy  $\bar{\sigma}(\pi^*, (\mathbf{u}_i^*)_{i \in \mathcal{I}})$  is a Markovian policy.*

PROOF. By Lemma 7, under the policy  $\bar{\sigma}(\pi^*, (\mathbf{u}_i^*)_{i \in \mathcal{I}})$ ,  $\pi^*$  is a stationary distribution, and  $(\mathbf{u}_i^*)_{i \in \mathcal{I}}$  are the corresponding transition frequencies. Recall the single-server  $\overline{\mathcal{P}}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$  problem

$$\begin{aligned} &\underset{\bar{N}, \bar{\sigma}, \pi}{\text{minimize}} && \bar{N} \\ &\text{subject to} && \mathbb{E}[h(\bar{\mathbf{K}}(\infty)) | \bar{\mathbf{K}}(\infty) \neq \mathbf{0}] \leq \epsilon, \\ &&& \bar{N} \cdot \bar{\lambda}_i = \lambda_i r, \quad \forall i \in \mathcal{I}. \end{aligned} \tag{3 revisited}$$

Observe that under the policy  $\bar{\sigma}(\pi^*, (\mathbf{u}_i^*)_{i \in \mathcal{I}})$ , the cost rate of resource contention is  $\mathbf{h}^T \pi \leq \epsilon(1 - \pi(\mathbf{0}))$ , the request rate of type  $i$  jobs is  $\bar{\lambda}_i = \mathbf{1}_o^T \mathbf{u}_i^* = \Phi^* \cdot \lambda_i$ , so

$$\begin{aligned} \mathbb{E}[h(\bar{\mathbf{K}}(\infty)) | \bar{\mathbf{K}}(\infty) \neq \mathbf{0}] &= \frac{\mathbf{h}^T \pi}{1 - \pi(\mathbf{0})} \leq \epsilon, \\ \bar{\lambda}_i &= \Phi^* \cdot \lambda_i, \quad \forall i \in \mathcal{I}. \end{aligned}$$

where we have used the fact that  $h(\mathbf{0}) = 0$  in the first equality. Therefore,  $(\Phi^*/r, \bar{\sigma}(\pi^*, (\mathbf{u}_i^*)_{i \in \mathcal{I}}), \pi)$  is a feasible solution to the single-server problem  $\overline{\mathcal{P}}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$ , achieving the objective value of  $r/\Phi^*$ , which is the optimal value because  $r/\Phi^* \leq \bar{N}^*$ .  $\square$

## D PERFORMANCE GUARANTEE OF JOIN-REQUESTING-SERVER WITH AN ESTIMATED MODEL

### D.1 Assumptions and result

In this section, we consider the performance of JRS when it is based on an estimated model. We will state the performance guarantee in terms of the estimation error, and give a proof sketch by pointing out which part of Theorem 3's proof needs to be changed accordingly.

Consider the setting where the maximal jobs on a server  $K_{\max}$ , the set of job phases  $\mathcal{I}$ , the cost rate function  $h(\cdot)$ , and the budget  $\epsilon$  are all known. However, we only have estimations of the jobs' arrival rates, internal transition rates, and departure rates.

Specifically, for any  $i, i' \in \mathcal{I}$ , let  $\tilde{\mu}_{ii'}$  be the *true* rate of internal transition from phase  $i$  to phase  $i'$ ; let  $\tilde{\mu}_{i\perp}$  be the *true* departure rate of phase  $i$ ; let  $\tilde{\lambda}_i r$  be the *true* arrival rate of type  $i$  jobs. We let  $\mu_{ii'}$ 's,  $\mu_{i\perp}$ 's, and  $\lambda_i r$ 's be the *estimated* internal transition rates, departure rates, and arrival rates, respectively. We assume that there exists a small positive constant  $\delta$  that is independent of the scaling factor  $r$  such that the following assumptions hold.

**Assumption 1** ( $\delta$ -accurate estimation). For any  $i, i' \in \mathcal{I}$ ,

$$|\mu_{ii'} - \tilde{\mu}_{ii'}| \leq \delta, \quad (73)$$

$$|\mu_{i\perp} - \tilde{\mu}_{i\perp}| \leq \delta, \quad (74)$$

$$|\lambda_i - \tilde{\lambda}_i| \leq \delta. \quad (75)$$

**Assumption 2** (Scaling of the single-server objective value). Consider the single-server problem in (3). Let  $(\bar{N}, \bar{\sigma}, \pi)$  be a solution feasible to the single-server problem with the estimated parameters that are  $\delta$ -accurate, where  $\delta \in [0, \delta_{\max})$ , for some constant  $\delta_{\max} > 0$ . We assume that there exist constants  $0 < m_1 < m_2$  independent of  $\delta$  and  $r$  such that

$$m_1 r \leq \bar{N} \leq m_2 r.$$

**Assumption 3** ( $\delta$ -insensitivity of the optimal value). Consider the single-server problem in (3). Let the optimal value of the single-server problem with the estimated parameters be  $\bar{N}^*$ , where the estimated parameters are  $\delta$ -accurate; let the optimal value of the single-server problem with true parameters be  $\bar{N}_{\text{true}}^*$ . We assume that

$$\bar{N}^* \leq \bar{N}_{\text{true}}^* + \delta r.$$

We also assume that JRS can accurately simulate the virtual jobs.

**Assumption 4** (Accurate simulation). The virtual jobs simulated in JRS follow the true transition dynamics.

Given the above assumptions, we have the following proposition that states the optimality gap of JRS policy with estimated model parameters, which has a similar form as Theorem 3 for JRS under true model parameters.

**Proposition 1** (Optimality gap with model estimation). *Consider a stochastic bin-packing problem in service systems with time-varying job resource requirements. Let the infinite-server policy  $\sigma$  be JRS with subroutine  $\bar{\sigma}$ . Suppose  $\sigma$  is specified based on an estimated model satisfying Assumption 1, 2, 3, and 4, for  $\delta$  s.t.  $\delta \in [0, \delta_{\max})$ , where  $\delta_{\max}$  is some positive constant independent of  $r$ . Let  $\bar{N}$  be the objective value achieved by  $\bar{\sigma}$  in the single-server problem  $\bar{\mathcal{P}}((\lambda_i r)_{i \in \mathcal{I}}, \epsilon)$  with estimated parameters. Under  $\sigma$ , for any initial state, we have*

$$\left| \sum_{k \neq 0} \mathbb{E}[X_k] - \lceil \bar{N} \rceil \cdot \mathbb{P}(\bar{K} \neq \mathbf{0}) \right| = O(\sqrt{r}) + \delta \cdot O(r), \quad (76)$$

$$\left| \sum_{k \neq 0} h(k) \mathbb{E}[X_k] - \lceil \bar{N} \rceil \cdot \mathbb{E}[h(\bar{K})] \right| = O(\sqrt{r}) + \delta \cdot O(r), \quad (77)$$

where  $\bar{\mathbf{K}}$  denotes the steady-state configurations of the single-server system under  $\bar{\sigma}$  with the estimated parameters. If we let  $\bar{\sigma}$  be the optimal policy of the single-server problem with estimated parameters, for any initial state, we have

$$N(\sigma) \leq (1 + B\delta + O(r^{-0.5})) \cdot \bar{N}_{true}^* \quad (78)$$

$$C(\sigma) \leq (1 + B\delta + O(r^{-0.5})) \cdot \epsilon, \quad (79)$$

where  $B$  is a positive constant independent of  $r$ . In other words,  $\sigma$  is  $(1 + B\delta + O(r^{-0.5}), 1 + B\delta + O(r^{-0.5}))$ -optimal.

**Remark 3.** Note even if the single-server system with estimated parameters  $\mathbf{k}^0$ -irreducible under  $\bar{\sigma}$ , it is hard to guarantee that the original system is  $\mathbf{k}^0$ -irreducible due to the estimation errors. Consequently, the steady-state performance metrics  $N(\sigma)$  and  $C(\sigma)$  could depend on the system's initial state. Fortunately, our proof for Theorem 3 does not rely on the uniqueness of the stationary distribution, so we can adapt the proof to show the inequalities in Proposition 1 for any initial state.

**Remark 4.** Note that Assumption 4 ensures that the real jobs and virtual jobs on a server are indistinguishable, so that  $(\hat{\mathbf{K}}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  is still a Markov chain. Recall that  $\hat{\mathbf{K}}^\ell$  and  $\boldsymbol{\eta}^\ell$  denote the observed configuration and tokens for each normal server  $\ell$ , respectively (see Section 5.1). We suspect that Assumption 4 can be removed since our proof does not rely too much on  $(\hat{\mathbf{K}}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  being a Markov chain. However, removing the assumption requires a more careful and notationally heavy analysis. We argue that in practice, this assumption is not restrictive, as one can record the traces of the jobs that arrived in the past and resample virtual jobs from those traces.

## D.2 Lemmas

In this section, we give a proof sketch for Proposition 1 when the single-server system under  $\bar{\sigma}$  with estimated parameters is  $\mathbf{k}^0$ -irreducible. The argument of extending to the general case is essentially the same as Appendix B.3, so we omit it here.

On a high level, the proof of Proposition 1 is similar to that of Theorem 1. In particular, the key steps of the proof are to show that  $d(\mathbf{K}^{1:L}, \bar{\mathbf{K}}^{1:L}) = O(\sqrt{r})$  and  $\sum_{\ell=L+1}^{\infty} \sum_{i \in I} K_i^\ell = O(\sqrt{r})$  as  $r \rightarrow \infty$ , where  $L = \lceil \bar{N} \rceil$ , and  $\bar{\mathbf{K}}^{1:L}$  are i.i.d. copies of steady-state configurations of the single-server system under  $\bar{\sigma}$  with the estimated parameters, and  $\mathbf{K}$  is the steady-state configurations of the infinite server system under  $\sigma$ .

Proposition 1 is based on the three lemmas stated below. The proof of Proposition 1 using the three lemmas is essentially the same as the argument in Section 5.3.

**Lemma 8.** Under the conditions of Proposition 1 and the single-server system with estimated parameters under  $\bar{\sigma}$  being  $\mathbf{k}^0$ -irreducible, for any initial state, we have

$$d(\hat{\mathbf{K}}^{1:L}, \bar{\mathbf{K}}^{1:L}) = O(\sqrt{r}) + \delta \cdot O(r)$$

**Lemma 9.** Under the conditions of Proposition 1 and the single-server system with estimated parameters under  $\bar{\sigma}$  being  $\mathbf{k}^0$ -irreducible, for any initial state and  $i \in I$ , the steady-state expected number of virtual jobs of type  $i$  s.t.

$$\mathbb{E} \left[ \sum_{\ell=1}^L \zeta_i^\ell \right] = O(\sqrt{r}) + \delta \cdot O(r).$$

**Lemma 10.** Under the conditions of Proposition 1 and the single-server system with estimated parameters under  $\bar{\sigma}$  being  $\mathbf{k}^0$ -irreducible, for any initial state and  $i \in I$ , the steady-state expected number of type  $i$  jobs on backup servers s.t.

$$\mathbb{E} \left[ \sum_{\ell=L+1}^{\infty} K_i^\ell \right] = O(\sqrt{r}) + \delta \cdot O(r).$$

These three lemmas are analogous to Lemma 2, Lemma 3, and Lemma 4, respectively. In the rest of the section, we sketch the proofs for the three lemmas above, highlighting the difference from the proofs of their analogues.

### D.3 Proof sketch for Lemma 8

Recall from Section 5.4 that the proof of Lemma 2 is based on Stein's method, which compares the generator of the i.i.d. copies of the single-server system with the generator of the infinite-server system. To write down the generators with the estimated model, recall that in the single-server system, each transition can be represented by the diagram

$$\mathbf{k} \rightarrow \mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a},$$

where the arrow  $\mathbf{k} \rightarrow \mathbf{k}'$  denotes an internal transition or a departure if  $\mathbf{k} \neq \mathbf{k}'$ ; the arrow  $\mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a}$  denotes a job request that is made right after reaching  $\mathbf{k}'$ . For any  $\mathbf{k}$ , let  $E(\mathbf{k})$  be the set of  $(\mathbf{k}', \mathbf{a}) \in \mathcal{K}^2$  such that  $\mathbf{k}' + \mathbf{a} \in \mathcal{K}$ . We define two sets of transition rates as below: for any  $\mathbf{k} \in \mathcal{K}$  and  $(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k})$ ,

- Under the estimated parameters and the policy  $\bar{\sigma}$ , we let  $\gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  be the rate of the transition  $\mathbf{k} \rightarrow \mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a}$ , and let  $\gamma_{\mathbf{k}} \triangleq \sum_{\mathbf{k}',\mathbf{a}} \gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  be the total transition rate at configuration  $\mathbf{k}$ .
- Under the *true parameters* and the policy  $\bar{\sigma}$ , we let  $\tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  be the rate of the transition  $\mathbf{k} \rightarrow \mathbf{k}' \rightarrow \mathbf{k}' + \mathbf{a}$ , and let  $\tilde{\gamma}_{\mathbf{k}} \triangleq \sum_{\mathbf{k}',\mathbf{a}} \tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  be the total transition rate at configuration  $\mathbf{k}$ .

Let  $\bar{G}$  be the generator of  $L = \lceil \bar{N} \rceil$  i.i.d. copies of single-server systems under the estimated parameters. For any  $g: \mathcal{K}^L \rightarrow \mathbb{R}$ , we have

$$\bar{G}g(\mathbf{k}^{1:L}) = \sum_{\ell=1}^L \sum_{\mathbf{k}',\mathbf{a}} \gamma_{\mathbf{k}^\ell,(\mathbf{k}',\mathbf{a})} (g(\cdot, \mathbf{k}' + \mathbf{a}, \cdot) - g(\cdot, \mathbf{k}^\ell, \cdot)), \quad (80)$$

where  $\sum_{\mathbf{k}',\mathbf{a}}$  is a shorthand for  $\sum_{(\mathbf{k}',\mathbf{a}) \in E(\mathbf{k})}$ . Let  $\hat{G}$  be the generator of  $(\hat{\mathbf{K}}^{1:L}(t), \boldsymbol{\eta}^{1:L}(t))$  for the infinite-server system. For any function  $g: \mathcal{K}^L \rightarrow \mathbb{R}$  and  $\psi(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L}) = g(\mathbf{k}^{1:L} + \boldsymbol{\eta}^{1:L})$ , we have

$$\begin{aligned} \hat{G}\psi(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L}) &= \sum_{\ell=1}^L \sum_{\mathbf{k}',\mathbf{a}} \tilde{\gamma}_{\mathbf{k}^\ell,(\mathbf{k}',\mathbf{a})} (g(\cdot, \mathbf{k}' + \mathbf{a}, \cdot) - g(\cdot, \mathbf{k}^\ell, \cdot)) \mathbb{1}_{\{\eta^\ell=0\}} \\ &\quad + \sum_{\ell=1}^L \sum_{\mathbf{k}',\mathbf{a}} \tilde{\gamma}_{\mathbf{k}^\ell,(\mathbf{k}',\mathbf{a})} (g(\cdot, \mathbf{k}' + \boldsymbol{\eta}^\ell, \cdot) - g(\cdot, \mathbf{k}^\ell + \boldsymbol{\eta}^\ell, \cdot)) \mathbb{1}_{\{\eta^\ell \neq 0\}}. \end{aligned} \quad (81)$$

To prove Lemma 8, we need to show that for any  $f \in \text{Lip}(1)$ ,  $\mathbf{k}^{1:L} \in \mathcal{K}^L$ , and  $\boldsymbol{\eta}^{1:L} \in \mathcal{K}^L$ ,

$$\bar{G}g_f(\mathbf{k}^{1:L} + \boldsymbol{\eta}^{1:L}) - \hat{G}\psi_f(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L}) = O(\sqrt{r}) + \delta \cdot O(r), \quad (82)$$

where  $g_f$  is the solution to

$$\mathbb{E} \left[ f(\bar{\mathbf{K}}^{1:L}) \right] - f(\mathbf{k}^{1:L}) = \bar{G}g_f(\mathbf{k}^{1:L}), \quad (83)$$

and  $\psi_f(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L}) = g_f(\mathbf{k}^{1:L} + \boldsymbol{\eta}^{1:L})$ . Same as the proof of Lemma 2, we prove (82) in two steps: the generator comparison step, and the Stein factor bound step.

In the generator comparison step, we observe that the formula of  $\bar{G}g(\mathbf{k}^{1:L})$  and  $\hat{G}\psi(\mathbf{k}^{1:L}, \boldsymbol{\eta}^{1:L})$  in (80) and (81) look almost the same as (15) and (17), except that the rates in (81) are  $\tilde{\gamma}_{\mathbf{k}^\ell,(\mathbf{k}',\mathbf{a})}$  instead of  $\gamma_{\mathbf{k}^\ell,(\mathbf{k}',\mathbf{a})}$ . As a result, after carrying out similar calculations as in the poof of Lemma 2, we get an extra error term involving  $\gamma_{\mathbf{k}^\ell,(\mathbf{k}',\mathbf{a})} - \tilde{\gamma}_{\mathbf{k}^\ell,(\mathbf{k}',\mathbf{a})}$ , which can be bounded using the lemma below. This error term results in the  $\delta \cdot O(r)$  in (82).

**Lemma 11.** *Under Assumption 1, for any  $\mathbf{k} \in \mathcal{K}$  and  $(\mathbf{k}', \mathbf{a}) \in E(\mathbf{k})$ , we have*

$$|\tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})} - \gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}| \leq K_{\max}\delta. \quad (84)$$

PROOF. When  $\mathbf{k} = \mathbf{k}'$ ,  $\gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  and  $\tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  are both the rate of adding  $\mathbf{a}$  jobs via a proactive request under the policy  $\bar{\sigma}$ , so they are identical.

When  $\mathbf{k} \neq \mathbf{k}'$  and  $\mathbf{a} = \mathbf{0}$ ,  $\gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  is equal to the rate of going from  $\mathbf{k}$  to  $\mathbf{k}'$  via an internal transition or departure under the estimated job model, while  $\tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  is under the true job model. Because there are at most  $K_{\max}$  jobs, and by our assumption the estimation error of each job's transition rates are bounded by  $\delta$ , thus  $\gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  and  $\tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  differ by at most  $K_{\max}\delta$ .

When  $\mathbf{k} \neq \mathbf{k}'$  and  $\mathbf{a} \neq \mathbf{0}$ ,  $\gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  is equal to the rate of going from  $\mathbf{k}$  to  $\mathbf{k}'$  via an internal transition or departure, multiplied by the probability of adding  $\mathbf{a}$  jobs via a reactive request, under the estimated job model and the policy  $\bar{\sigma}$ ;  $\tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  is under the true job model instead of the estimated job model, but uses the same policy. Because the rate of going from  $\mathbf{k}$  to  $\mathbf{k}'$  differs by at most  $K_{\max}\delta$  under two different job models, and the probability of adding  $\mathbf{a}$  jobs after going from  $\mathbf{k}$  to  $\mathbf{k}'$  is the same under the same policy, thus  $\gamma_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  and  $\tilde{\gamma}_{\mathbf{k},(\mathbf{k}',\mathbf{a})}$  differ by at most  $K_{\max}\delta$ .  $\square$

In the Stein factor bound step, we need to show that

$$\sup_{\mathbf{k}, \mathbf{k}' \in \mathcal{K}} |g_f(\cdot, \mathbf{k}', \cdot) - g_f(\cdot, \mathbf{k}, \cdot)| = O(1). \quad (85)$$

This involves analyzing the i.i.d. copies of the single-server system, and the fact that the single-server system with estimated parameters is  $\mathbf{k}^0$ -irreducible under  $\bar{\sigma}$ . This part of the proof is identical to the corresponding part of the proof of Lemma 2.

#### D.4 Proof sketch for Lemma 9 and Lemma 10

The proof of Lemma 9 and Lemma 10 has a similar structure as the proof of Lemma 3 and Lemma 4. In the first step, we use Little's law to bound expectations of the number of type  $i$  virtual jobs,  $V_i$ , and the number of type  $i$  jobs on backup servers,  $Y_i$ . We have two equations almost the same as (36) and (37) except that the rates  $\gamma_{\tilde{\mathbf{K}}^\ell,(\mathbf{k}',\mathbf{a})}$  and  $\lambda_i r$  are replaced by  $\tilde{\gamma}_{\tilde{\mathbf{K}}^\ell,(\mathbf{k}',\mathbf{a})}$  and  $\tilde{\lambda}_i r$ :

$$\mathbb{E}[V_i] \leq t_{\max} \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{\gamma}_{\tilde{\mathbf{K}}^\ell,(\mathbf{k}',\mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right], \quad (86)$$

$$\mathbb{E}[Y_i] \leq t_{\max} \mathbb{E} \left[ \tilde{\lambda}_i r \cdot dy_i \right], \quad (87)$$

where  $t_{\max}$  is the maximal expected service time of any type of virtual job or real job. Because  $t_{\max} = O(1)$ , it suffices to bound  $\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{\gamma}_{\tilde{\mathbf{K}}^\ell,(\mathbf{k}',\mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right]$  and  $\mathbb{E} \left[ \tilde{\lambda}_i r \cdot dy_i \right]$ .

In the second step, we utilize the fact that the two Lyapunov functions  $g(z_i) = z_i$  and  $g(z_i) = z_i^2$  have zero drift in steady-state, where  $z_i$  is the total number of type  $i$  tokens. We get two equalities similar to (38) and (41):

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{\gamma}_{\tilde{\mathbf{K}}^\ell,(\mathbf{k}',\mathbf{a})} (a_i - dv_i) \mathbb{1}_{\{\eta^\ell=0\}} + \tilde{\lambda}_i r (-1 + dy_i) \right] = 0. \quad (88)$$

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{\gamma}_{\tilde{\mathbf{K}}^\ell,(\mathbf{k}',\mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right] \quad (89)$$

$$= \frac{1}{\eta_{\max}} \cdot \mathbb{E} \left[ \left( \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{\gamma}_{\tilde{\mathbf{K}}^\ell,(\mathbf{k}',\mathbf{a})} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \tilde{\lambda}_i r \right) \cdot Z_i \right] \quad (90)$$



$$+ \frac{1}{2\eta_{\max}} \cdot \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} (a_i^2 - (dv_i)^2) \mathbb{1}_{\{\eta^\ell=0\}} + \tilde{\lambda}_i r \cdot (1 - (dy_i)^2) \right]. \quad (91)$$

In the third step, we use the above two equalities to bound  $\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right]$  and  $\mathbb{E} \left[ \tilde{\lambda}_i r \cdot dy_i \right]$ . We first use the equality in (89) to (91) to bound  $\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right]$ . Following the same argument as the proof of Lemma 3 and Lemma 4 until (47), we can show that

$$(91) \leq O(\sqrt{r}),$$

$$\begin{aligned} (90) &\leq \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \mathbb{1}_{\{\eta^\ell=0\}} - \tilde{\lambda}_i r \right\| \\ &\leq \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i - \tilde{\lambda}_i r \right\| + O(\sqrt{r}) \\ &\leq \mathbb{E} \left\| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i - \tilde{\lambda}_i r \right\| + O(\sqrt{r}) + \delta \cdot O(r), \end{aligned} \quad (92)$$

where to get (92), we apply Lemma 8 to replace  $\widehat{\mathbf{K}}^\ell$  with  $\overline{\mathbf{K}}^\ell$ , which causes an  $O(\sqrt{r}) + \delta \cdot O(r)$  error. Next, we show that

$$\mathbb{E} \left\| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i - \tilde{\lambda}_i r \right\| = O(\sqrt{r}) + \delta \cdot O(r). \quad (93)$$

By Assumption 1 and Lemma 11, we have

$$|\tilde{\lambda}_i r - \lambda_i r| \leq \delta r, \quad (94)$$

$$\left| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i - \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} Y_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \right| \leq \delta \cdot O(r). \quad (95)$$

These two bounds allow us to replace  $\tilde{Y}_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})}$  and  $\tilde{\lambda}_i r$  on the LHS of (93) with  $Y_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})}$  and  $\lambda_i r$  at the cost of introducing  $\delta \cdot O(r)$  error. Moreover, because  $\{\overline{\mathbf{K}}^\ell\}_{\ell=1, \dots, L}$  are i.i.d.,  $\sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} Y_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i$  concentrates around its mean with  $O(\sqrt{r})$  error, where the mean can be shown to be

$$\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} Y_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i \right] = \bar{\lambda}_i L = \lambda_i r + O(1).$$

Note that the first equality in (96) holds because for each  $\ell$ ,  $\mathbb{E}[\sum_{\mathbf{k}', \mathbf{a}} Y_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i]$  is equal to  $\bar{\lambda}_i$ , i.e., the long-run average rate of requesting type  $i$  jobs on a single-server system with estimated parameters under  $\bar{\sigma}$ ; the second equality in (96) is because  $L = \lceil \bar{N} \rceil$ , and  $\bar{N} \bar{\lambda}_i = \lambda_i r$ . As a result,

$$\mathbb{E} \left\| \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} Y_{\overline{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} a_i - \lambda_i r \right\| \leq O(\sqrt{r}). \quad (96)$$

Combining (94) to (96), we get (93). Therefore,

$$\mathbb{E}[V_i] \leq t_{\max} \mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', \mathbf{a}} \tilde{Y}_{\widehat{\mathbf{K}}^\ell, (\mathbf{k}', \mathbf{a})} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right]$$

$$\begin{aligned}
&\leq O(1) \cdot ((90) + (91)) \\
&\leq O(\sqrt{r}) + \delta \cdot O(r),
\end{aligned}$$

which completes the proof of Lemma 9.

Finally, it is straightforward to show that  $\mathbb{E} \left[ \tilde{\lambda}_i r \cdot dy_i \right] = O(\sqrt{r}) + \delta \cdot O(r)$  using (88) and the bound on  $\mathbb{E} \left[ \sum_{\ell=1}^L \sum_{\mathbf{k}', a} \tilde{Y}_{\tilde{K}^\ell, (\mathbf{k}', a)} dv_i \mathbb{1}_{\{\eta^\ell=0\}} \right]$ . Therefore,  $\mathbb{E}[Y_i] \leq t_{\max} \mathbb{E} \left[ \tilde{\lambda}_i r \cdot dy_i \right] = O(\sqrt{r}) + \delta \cdot O(r)$ . This proves Lemma 10.

Received August 2023; revised October 2023; accepted October 2023