



ISE

Industrial and
Systems Engineering

A service system with packing constraints: Greedy
randomized algorithm achieving sublinear in scale
optimality gap

ALEXANDER L. STOLYAR¹ AND YUAN ZHONG²

¹Lehigh University ²Columbia University

ISE Technical Report 15T-019



LEHIGH
UNIVERSITY.

A service system with packing constraints: Greedy randomized algorithm achieving sublinear in scale optimality gap

Alexander L. Stolyar
Lehigh University
200 W. Packer Ave., Mohler 484
Bethlehem, PA 18015
stolyar@lehigh.edu

Yuan Zhong
Columbia University
500 W. 120 St, Mudd 344
New York, NY 10027
yz2561@columbia.edu

November 10, 2015

Abstract

A service system with multiple types of arriving customers is considered. There is an infinite number of homogeneous servers. Multiple customers can be placed for simultaneous service into one server, subject to general *packing constraints*. The service times of different customers are independent, even if they are served simultaneously by the same server; the service time distribution depends on the customer type. Each new arriving customer is placed for service immediately, either into an occupied server, i.e., one already serving other customers, as long as packing constraints are not violated, or into an empty server. After service completion, each customer leaves its server and the system. The basic objective is to minimize the number of occupied servers in steady state.

We study a *Greedy-Random* (GRAND) placement (packing) algorithm, introduced in [23]. This is a simple online algorithm, which places each arriving customer uniformly at random into either one of the already occupied servers that can still fit the customer, or one of the so-called *zero-servers*, which are empty servers designated to be available to new arrivals. In [23], a version of the algorithm, labeled GRAND(aZ), was considered, where the number of zero servers is aZ , with Z being the current total number of customers in the system, and $a > 0$ being an algorithm parameter. GRAND(aZ) was shown in [23] to be asymptotically optimal in the following sense: (a) the steady-state optimality gap grows linearly in the system scale r (the mean total number of customers in service), i.e. as $c(a)r$ for some $c(a) > 0$; and (b) $c(a) \rightarrow 0$ as $a \rightarrow 0$.

In this paper, we consider the GRAND(Z^p) algorithm, in which the number of zero-servers is Z^p , where $p \in (1 - 1/(8\kappa), 1)$ is an algorithm parameter, and $(\kappa - 1)$ is the maximum possible number of customers that a server can fit. We prove the asymptotic optimality of GRAND(Z^p) in the sense that the steady-state optimality gap is $o(r)$, sublinear in the system scale. This is a stronger form of asymptotic optimality than that of GRAND(aZ).

1 Introduction

Efficient resource allocation in modern cloud computing systems poses many interesting and challenging new problems; see e.g., [10] for an overview. One such problem is that of efficient real-time assignment (or, packing) of virtual machines to physical machines in a cloud data center, where a primary objective is to minimize the number of physical machines being used. This leads to stochastic dynamic bin packing models, where, in contrast to many classical bin packing models, “items” (virtual machines) being placed into “bins” (physical machines) do not stay in the system forever, but leave after a random “service time.” For this reason, such models are naturally viewed as service (or queueing) systems, with “items” and “bins” being viewed as customers and servers, respectively.

In this paper, we consider a general service system model that has also been studied in [21–23]. There is a finite number of customer types, and the number of available servers is infinite. Customers arrive to the system over time, and multiple customers can be placed for simultaneous service (or *fit*) into the same server, subject to *packing constraints*. We consider *monotone* packing, a general class of packing constraints, where we only impose the following natural and non-restrictive *monotonicity* condition: if a certain set of customers can fit into a server, then a subset can fit as well. The servers are *homogeneous* in that they all have the same packing constraints. The service times of different customers are independent, even if they are served simultaneously by the same server, and the service time distribution depends only on the customer type. Each new arriving customer is placed for service immediately, either into an occupied server, i.e., one that is already serving other customers, as long as packing constraints are not violated, or into an empty server. After service completion, each customer leaves its server and the system – as we mentioned, this is what distinguishes our model from many classical bin packing models (see, e.g., [1, 7]). As in [23], we make Markov assumptions, where customers of each type arrive as an independent Poisson process, and service time distributions are all exponential.

We are interested in designing customer placement (packing) algorithms that minimize the total number of occupied servers in steady state. In addition, given the scale at which modern cloud data centers operate, it is highly desirable that a placement algorithm is *online*, i.e., it makes decisions based on the current system state only, and *parsimonious*, i.e., it requires only minimal knowledge of system structure and state information.

We study a *Greedy-Random* (GRAND) placement algorithm, introduced in [23]. This is a very parsimonious online algorithm, which places each arriving customer uniformly at random into either one of the already occupied servers (subject to packing constraints), or one of the so-called *zero-servers*, which are empty servers designated to be available to new arrivals. In [23], a version of the algorithm, which we call $\text{GRAND}(aZ)$, was considered, where the number of zero servers depend on the current system state as aZ , with Z being the current total number of customers in the system, and $a > 0$ being an algorithm parameter. $\text{GRAND}(aZ)$ was shown in [23] to be asymptotically optimal, in the sense of the following two properties: (a) for each $a > 0$, the steady-state optimality gap is of the form $c^r(a)r$, where r is the system *scale*, defined to be the expectation of Z in steady state, with $c^r(a) \rightarrow c(a) > 0$ as $r \rightarrow \infty$; and (b) $c(a) \rightarrow 0$ as $a \rightarrow 0$. In other words, under $\text{GRAND}(aZ)$ the optimality gap grows linearly as the system scale r , with the linear factor going to 0 as the algorithm parameter $a \rightarrow 0$.

The focus of this paper is the $\text{GRAND}(Z^p)$ algorithm, in which the number of zero-servers is Z^p , where $p \in (1 - \frac{1}{8\kappa}, 1)$ is an algorithm parameter, and $(\kappa - 1)$ is the maximum possible number of customers that a server can fit. Our **main result** is the asymptotic optimality of $\text{GRAND}(Z^p)$, in the sense that the steady-state optimality gap is $o(r)$, i.e. it is sublinear in the system scale. This is a stronger form of asymptotic optimality than that of $\text{GRAND}(aZ)$, because $\text{GRAND}(Z^p)$ achieves the sublinear gap simply as $r \rightarrow \infty$, without having to take an additional limit on any algorithm parameter. This is in contrast to the case of $\text{GRAND}(aZ)$, which achieves asymptotic optimality only by first taking the limit $r \rightarrow \infty$, and then the limit $a \rightarrow 0$ on the algorithm parameter a .

Let us provide some remarks on the reasons why this stronger form of the asymptotic optimality of $\text{GRAND}(Z^p)$ is, on the one hand, natural to expect, and, on the other hand, substantially more difficult to rigorously prove.

Informally speaking, $\text{GRAND}(Z^p)$ can be viewed as $\text{GRAND}(aZ)$, where the parameter a , instead of being fixed, is replaced by a variable that decreases to zero with increasing system scale; namely, $a = Z^p/Z = Z^{p-1}$. However, the asymptotic optimality of $\text{GRAND}(aZ)$ does *not*, of course, imply the (stronger form of) asymptotic optimality of $\text{GRAND}(Z^p)$, because the system scale $r \approx Z$, and, therefore, the “parameter” $a = Z^{p-1} \approx r^{p-1}$ depends on and changes with the scale r . The key technical difficulty is that the analysis of $\text{GRAND}(Z^p)$ *cannot* be reduced to the analysis of system fluid limits and/or local fluid limits. See Section 6.1 for a more detailed discussion.

$\text{GRAND}(Z^p)$, just like $\text{GRAND}(aZ)$, is extremely simple and easy to implement. It does not need to keep track of the exact configurations of the servers, and the only information required at any time is which customer types a given server can still accept for service. As a result, the algorithm only needs to maintain a very small number of variables. Furthermore, the algorithm does not use knowledge of the customer arrival rates or expected service times. We refer the readers to [23] for a more detailed discussion of these attractive implementational features.

1.1 Brief Literature Review

Our work is related to several lines of previous research. First, our work is related to the extensive literature on classical bin packing problems, where items of various sizes arrive to the system, and need to be placed in finite-size bins, according to an online algorithm. Once placed, items never leave or move between bins. The *worst-case analysis* of such problems considers all possible instances of item sizes and sequencings, and aims to develop simple algorithms with performance guarantee over all problem instances. See e.g., [4] for a recent, extensive survey. The *stochastic analysis* assumes that item sizes are given according to a probability distribution, and the typical objective is to minimize the expected number of occupied bins. For an overview of results, see e.g., [7] and references therein. A recent paper [12] establishes improved results for the classical stochastic setting, and contains some heuristics and simulations for the case with item departures, which is a special case of our model.

Much research on classical bin packing concerns the one-dimensional case, in which both item and bin sizes are scalars. It is possible to generalize one-dimensional bin packing to higher dimensions in multiple ways. For example, in vector packing problems [3,21], item and bin sizes are vectors, and in box packing problems [1], items and bins are rectangles or hyper-rectangles. Let us also remark that the packing constraints in our model include vector packing and box packing as special cases. See e.g., [1] for an overview of multi-dimensional packing.

Motivated partly by applications to computer storage allocation, a *one-dimensional, dynamic bin packing* problem was introduced in [5], where, similar to the model of this paper, items leave the system after a finite service time. Paper [5] contains a worst-case analysis of the problem, so the techniques used are quite distinct from those in our paper. A recent review of results on worst-case dynamic bin packing can be found in [4].

Another related line of works considers bin packing *service* systems, which have one (see e.g. [6,8]) or several servers (see e.g., [11,13,15,16]). In these systems, random-size items (or customers) arrive over time, and get placed into servers for processing, subject to packing constraints. Customers waiting for service are queued, and a typical problem is to determine the maximum throughput and/or minimum queueing delay under a packing algorithm. Our model is similar to these systems, since they model customer departures, but is also different, mainly because our system has an infinite number of servers, so that there are no queues or problem of stability. Like our work, recent papers on bin packing service systems with multiple servers [11, 13, 15, 16] are also motivated by real-time VM placement problems.

As mentioned in the introduction, the model in this paper is the same as that studied in [21–23]. Papers [21,22] introduce and study different classes of *Greedy* algorithms, and prove their asymptotic optimality, as the system scale grows to infinity. A Greedy algorithm does not use the knowledge of the customer arrival rates or mean service times, and makes placement decisions based on the current system state only. However,

unlike GRAND, it does need to keep track of the numbers of servers in different *packing configurations*. The number of possible configurations can be prohibitively large in many practical scenarios, a feature that may limit the implementability of Greedy.

The relation of our main results to those for GRAND(aZ) in [23] has already been discussed in much detail. Also related to the GRAND algorithms is a recent paper [24], which generalizes GRAND(aZ) and its asymptotic optimality to *heterogeneous* systems with multiple server types, in which packing constraints depend on the server type; it also contains results for heterogeneous systems with *finite pools* of servers of each type. Finally, a randomized version of the Best-Fit algorithm was studied in [9], which considers the model of this paper with specialized packing constraints, and was proved to be asymptotically optimal, using techniques similar to those in [23].

1.2 Organization

The rest of the paper is organized as follows. In Subsection 1.3, we introduce basic notation and conventions that will be used throughout the paper. The model, the GRAND(Z^p) algorithm and the asymptotic regime are formally defined in Section 2. We state Theorem 4, the main result on the asymptotic optimality of GRAND(Z^p), in Section 3. The basic system dynamics, together with additional notation and terminology, are described in Section 4. In Section 5, we collect some results and observations obtained in [23], which are needed for our proofs. The proof of Theorem 4 is in Section 6, with the proofs of some auxiliary results given in the Appendix. The paper is concluded in Section 7 with some discussion and suggestions for future work.

1.3 Basic Notation and Conventions

Sets of real and non-negative real numbers are denoted by \mathbb{R} and \mathbb{R}_+ , respectively. Similarly, sets of integers and non-negative integers are denoted by \mathbb{Z} and \mathbb{Z}_+ , respectively. For $\xi \in \mathbb{R}$, $\lceil \xi \rceil$ denotes the smallest integer greater than or equal to ξ , and $\lfloor \xi \rfloor$ denotes the largest integer smaller than or equal to ξ . For two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we use $\mathbf{x} \cdot \mathbf{y}$ to denote their scalar (dot) product; i.e., $\mathbf{x} \cdot \mathbf{y} = \sum_i x_i y_i$. The standard Euclidean norm of a vector $\mathbf{x} \in \mathbb{R}^n$ is denoted by $\|\mathbf{x}\| = \sqrt{\mathbf{x} \cdot \mathbf{x}}$, and the distance from vector \mathbf{x} to a set U in a Euclidean space is denoted by $d(\mathbf{x}, U) = \inf_{\mathbf{u} \in U} \|\mathbf{x} - \mathbf{u}\|$. $\mathbf{x} \rightarrow \mathbf{u} \in \mathbb{R}^n$ means ordinary convergence in \mathbb{R}^n , and $\mathbf{x} \rightarrow U \subseteq \mathbb{R}^n$ means $d(\mathbf{x}, U) \rightarrow 0$. For $\mathbf{x} \in \mathbb{R}^n$, we also use $\|\mathbf{x}\|_1$ to denote the 1-norm of \mathbf{x} , defined to be $\|\mathbf{x}\|_1 = \sum_i |x_i|$. \mathbf{e}_i is the i -th coordinate unit vector in \mathbb{R}^n . Symbol \implies denotes convergence in distribution of random variables taking values in space \mathbb{R}^n equipped with the Borel σ -algebra. Symbol $\stackrel{d}{=}$ means *equal in distribution*. The abbreviation *w.p.1* means *with probability 1*. The abbreviation RHS (LHS, respectively) means *right-hand side* (*left-hand side*, respectively). In addition, the abbreviation *w.r.t* means *with respect to*, and abbreviation *u.o.c.* means *uniform on compact sets*. We often write $x(\cdot)$ to mean the function (or random process) $\{x(t), t \geq 0\}$, and we write $\{x_{\mathbf{k}}\}$ to mean the vector $\{x_{\mathbf{k}}, \mathbf{k} \in \mathcal{K}\}$, where the set of indices \mathcal{K} is determined by the context. For a function (or random process) $x(\cdot)$, we use $x(t+)$ to denote its right limit at time t , i.e., $x(t+) = \lim_{s \downarrow t} x(s)$, and use $x(t-)$ to denote its left limit at time t , i.e., $x(t-) = \lim_{s \uparrow t} x(s)$, whenever these limits exist. For a differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$, its gradient at $\mathbf{x} \in \mathbb{R}^n$ is denoted by $\nabla f(\mathbf{x}) = \left\{ \frac{\partial}{\partial x_i} f(\mathbf{x}), i = 1, \dots, n \right\}$. Indicator function $I\{A\}$ for a condition A is equal to 1 if A holds and 0 otherwise. The cardinality of a finite set \mathcal{N} is $|\mathcal{N}|$. Notation \doteq means *is defined to be*. In this paper, we use bold letters to represent vectors, and plain letters to represent scalars.

2 System Model

2.1 Infinite Server System with Packing Constraints

The model that we study in this paper is the same as in [21–23]. We consider a service system with an infinite number of servers. There are I types of arriving customers, indexed by $i \in \{1, 2, \dots, I\} \doteq \mathcal{I}$. For each $i \in \mathcal{I}$, customers of type i arrive as an independent Poisson process of rate $\Lambda_i > 0$, and have service times that are exponentially distributed with mean $1/\mu_i$. The service times of all customers are mutually independent, and independent of the arrival processes. Each customer is placed into one of the servers for processing, immediately upon arrival, and departs the system after the service completes. Multiple customers can occupy the same server simultaneously, subject to the so-called *monotone* packing constraint, which we now describe.

Definition 1 (Monotone packing). *A packing constraint is characterized by a finite set $\bar{\mathcal{K}}$, the set of feasible server configurations. A vector $\mathbf{k} = \{k_i, i \in \mathcal{I}\} \in \bar{\mathcal{K}}$ is a (feasible) server configuration if (a) for each $i \in \mathcal{I}$, $k_i \in \mathbb{Z}_+$; and (b) a server can simultaneously process k_1 customers of type 1, k_2 customers of type 2, \dots , and k_I customers of type I . The packing constraint $\bar{\mathcal{K}}$ is called monotone, if the following condition holds: whenever $\mathbf{k} \in \bar{\mathcal{K}}$ and $\mathbf{k}' \leq \mathbf{k}$ componentwise, then $\mathbf{k}' \in \bar{\mathcal{K}}$ as well.*

Without loss of generality, we assume that $\mathbf{e}_i \in \bar{\mathcal{K}}$ for all $i \in \mathcal{I}$, where \mathbf{e}_i is the i -th coordinate unit vector, so that customers of all types can be processed. By definition, the component-wise zero vector $\mathbf{0}$ belongs to $\bar{\mathcal{K}}$ – this is the configuration of an empty server. We denote by $\mathcal{K} = \bar{\mathcal{K}} \setminus \{\mathbf{0}\}$ the set of configurations that *do not* include the zero configuration. As discussed in [21, 22], monotone packing includes important special packing constraints such as vector packing.

An important assumption of the model is that simultaneous service does *not* affect the service time distributions of individual customers. In other words, the service time of a customer is unaffected by whether or not there are other customers served simultaneously by the same server.

We now define the system state. For each $\mathbf{k} \in \mathcal{K}$, let $X_{\mathbf{k}}(t)$ denote the number of servers with configuration \mathbf{k} . Then, the vector $\mathbf{X}(t) = \{X_{\mathbf{k}}(t), \mathbf{k} \in \mathcal{K}\}$ is the *system state* at time t , and we often write $\mathbf{X} = \{X_{\mathbf{k}}, \mathbf{k} \in \mathcal{K}\}$ for a generic state. Note that the system state does not include the number of empty servers, which would always be infinite.

2.2 GRAND and GRAND(Z^p) Algorithms

In general, a *placement (or packing) algorithm* determines the servers into which arriving customers are placed dynamically over time. In this paper, we are interested in *online* placement algorithms, which make placement decisions based only on the current system state \mathbf{X} . Thus, from now on, we will use “placement algorithms” to mean online ones. Our primary objective is the design of placement algorithms that minimize the total number of occupied servers $\sum_{\mathbf{k} \in \mathcal{K}} X_{\mathbf{k}}$ in the stationary regime.

Under any well-defined placement algorithm, the process $\{\mathbf{X}(t), t \geq 0\}$ is a continuous-time Markov chain with a countable state space, which is irreducible and positive recurrent. The positive recurrence of $\{\mathbf{X}(t), t \geq 0\}$ can be argued as follows. First, for each $i \in \mathcal{I}$ and each $t \geq 0$, we let

$$Y_i(t) = \sum_{\mathbf{k} \in \mathcal{K}} k_i X_{\mathbf{k}}(t) \tag{1}$$

be the number of type- i customers in the system at time t , and let

$$Z(t) = \sum_i Y_i(t) \tag{2}$$

be the total number of customers in the system at time t . For each $i \in \mathcal{I}$, $\{Y_i(t) : t \geq 0\}$ describes exactly the dynamics of an independent $M/M/\infty$ queueing system with arrival rate Λ_i and service rate μ_i , regardless of the placement algorithm, and has a unique stationary distribution. Let us denote by $Y_i(\infty)$ the random value of $Y_i(t)$ in the stationary regime. Then, $Y_i(\infty)$ is a Poisson random variable with mean Λ_i/μ_i for each $i \in \mathcal{I}$. Similarly, the process $\{Z(t) : t \geq 0\}$ also has a unique stationary distribution, and if we let $Z(\infty)$ be the random value of $Z(t)$ in the stationary regime, then $Z(\infty)$ is a Poisson random variable with mean $\sum_{i \in \mathcal{I}} \Lambda_i/\mu_i$. Thus, the positive recurrence of the Markov chain $\{\mathbf{X}(t), t \geq 0\}$ follows from the facts that $\sum_{\mathbf{k} \in \mathcal{K}} X_{\mathbf{k}}(t) \leq \sum_i Y_i(t) = Z(t)$ for all $t \geq 0$, and that the process $\{Z(t) : t \geq 0\}$ is positive recurrent. Consequently, the process $\{\mathbf{X}(t), t \geq 0\}$ has a unique stationary distribution, and we let $\mathbf{X}(\infty) = \{X_{\mathbf{k}}(\infty), \mathbf{k} \in \mathcal{K}\}$ be the random system state $\mathbf{X}(t)$ in the stationary regime.

In [23], we introduced a broad class of placement algorithms, called the Greedy-Random (GRAND) algorithms. For completeness, we include the definition here.

Definition 2 (Greedy-Random (GRAND) algorithm). *At any given time t , there is a designated finite set of $X_{\mathbf{0}}(t)$ empty servers, called zero-servers, where $X_{\mathbf{0}}(t) = f(\mathbf{X}(t))$ is a given fixed function of the system state $\mathbf{X}(t)$. Suppose that a type- i customer arrives at time t . Then, the customer is placed into a server chosen uniformly at random among the zero-servers and the occupied servers where the customer can fit. In other words, the total number of servers available to a type- i arrival at time t is*

$$X_{(i)}(t) \doteq X_{\mathbf{0}}(t) + \sum_{\mathbf{k} \in \mathcal{K}: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} X_{\mathbf{k}}(t).$$

If $X_{(i)}(t) = 0$, then the customer is placed into an empty server. Furthermore, it is important that immediately after any arrival or departure at time t , the value of $X_{\mathbf{0}}(t+)$ is reset (if necessary) to $f(\mathbf{X}(t+))$.

Note that the number $X_{\mathbf{0}} = f(\mathbf{X})$ of zero-servers is finite at all times, even though there is always an infinite number of empty servers available. Also, recall that a system state $\mathbf{X} = \{X_{\mathbf{k}}, \mathbf{k} \in \mathcal{K}\}$ only includes the quantities of servers in non-zero configurations, so it does not include $X_{\mathbf{0}}$. Clearly, for a given function $f(\cdot)$, GRAND is a well-defined placement algorithm.

The focus of this paper is the following specialization of the general GRAND algorithm.

Definition 3 (GRAND(Z^p) algorithm). *A special case of the GRAND algorithm, with the number of zero-servers depending only on the total number of customers as*

$$X_{\mathbf{0}}(t) = \lceil Z(t)^p \rceil \quad \forall t \geq 0,$$

where $p \in (0, 1)$ is a parameter, is called a GRAND(Z^p) algorithm.

In [23], we introduced a different specialization of the general GRAND algorithm, namely the GRAND(aZ) algorithm, where the number of zero-servers $X_{\mathbf{0}}$ depends on Z , the total number of customers, as $X_{\mathbf{0}} = \lceil aZ \rceil$, $a > 0$. GRAND(Z^p) is similar to GRAND(aZ) in that under both algorithms, $X_{\mathbf{0}}$ only depends on the system state \mathbf{X} through Z , the total number of customers.

As we can see, GRAND are fairly “blind” algorithms in that when they place a customer, they do not prefer one configuration over another, as long as they can fit this additional customer. Similar to GRAND(aZ), GRAND(Z^p) can be efficiently implemented. At all times, the algorithm only needs to keep track of $I + 1$ variables to make placement decisions, with the variables being (a) Z , the total number of customers; and (b) for each $i \in \mathcal{I}$, $X_{(i)}$, the total number of servers that can fit an additional type- i customer, including zero-servers and occupied ones. In contrast to GRAND(Z^p) and GRAND(aZ), the Greedy algorithm [21] and the Greedy-with-Sublinear-safety-Stocks (GSS) algorithm [22] both need to keep track of $|\mathcal{K}|$ number of variables, which is often prohibitively large in real applications.

2.3 Asymptotic Regime

We are interested in the asymptotic performance properties of the GRAND(Z^p) algorithms in the scaling regime $r \rightarrow \infty$. More specifically, assume that $r \geq 1$, and r increases to infinity along a discrete sequence. Customer arrival rates scale linearly with r ; i.e., for each r , $\Lambda_i = \lambda_i r$, where λ_i are fixed positive parameters. Without loss of generality, we assume that $\sum_i \lambda_i / \mu_i = 1$, by suitably re-defining r if required. For each r , let $\mathbf{X}^r(\cdot)$ be the random process associated with the system parametrized by r , let $\mathbf{X}^r(\infty)$ be the (random) system state in the stationary regime, and let $X_0^r(t)$ be the number of zero-servers in the system at time t . For each $i \in \mathcal{I}$, let $Y_i^r(t) = \sum_{\mathbf{k} \in \mathcal{K}} k_i X_{\mathbf{k}}^r(t)$ be the total number of type- i customers at time t , and let $Y_i^r(\infty) = \sum_{\mathbf{k} \in \mathcal{K}} k_i X_{\mathbf{k}}^r(\infty)$. Similarly, let $Z^r(t) = \sum_{i \in \mathcal{I}} Y_i^r(t)$ be the total number of customers at time t and let $Z^r(\infty) = \sum_i Y_i^r(\infty)$. As explained in Section 2.2, for each $i \in \mathcal{I}$, $Y_i^r(\infty)$ is an independent Poisson random variable with mean $r \rho_i$, where $\rho_i \equiv \lambda_i / \mu_i$, and $Z^r(\infty)$ is a Poisson random variable with mean $r \sum_i \rho_i = r$.

The *fluid-scaled* processes are defined to be $\{\mathbf{x}^r(t) : t \geq 0\}$ for each r , where $\mathbf{x}^r(t) = \mathbf{X}^r(t)/r$. We also define $\mathbf{x}^r(\infty) = \mathbf{X}^r(\infty)/r$. For any r , $\mathbf{x}^r(t)$ takes values in the non-negative orthant $\mathbb{R}_+^{|\mathcal{K}|}$. Similarly, $y_i^r(t) \doteq Y_i^r(t)/r$, $z^r(t) \doteq Z^r(t)/r$, $x_0^r(t) \doteq X_0^r(t)/r$ and $x_{(i)}^r(t) \doteq X_{(i)}^r(t)/r$, for $t \in [0, \infty]$. Since $\sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}}^r(\infty) \leq \sum_i y_i^r(\infty) \leq z^r(\infty) = Z^r(\infty)/r$, the family of random variables $\sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}}^r(\infty)$ is uniformly integrable in r . This in particular implies that the sequence of distributions of $\mathbf{x}^r(\infty)$ is tight, so there always exists a limit $\mathbf{x}(\infty)$ in distribution, and $\mathbf{x}^r(\infty) \implies \mathbf{x}(\infty)$, along a subsequence of r .

Since $Y_i^r(\infty)$ is a Poisson random variable of mean $\rho_i r$, any weak limit point $y_i(\infty)$ of $\{y_i^r(\infty)\}_r$ must concentrate at the constant ρ_i . Thus, the limit (random) vector $\mathbf{x}(\infty)$ satisfies the following conservation laws:

$$\sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}}(\infty) \equiv y_i(\infty) = \rho_i, \quad \forall i, \quad (3)$$

which, in particular, implies that

$$z(\infty) \equiv \sum_i y_i(\infty) \equiv \sum_i \rho_i = 1. \quad (4)$$

Therefore, the values of $\mathbf{x}(\infty)$ are confined to the convex compact $(|\mathcal{K}| - I)$ -dimensional polyhedron

$$\mathcal{X} \doteq \left\{ \mathbf{x} \in \mathbb{R}_+^{|\mathcal{K}|} \mid \sum_{\mathbf{k}} k_i x_{\mathbf{k}} = \rho_i, \quad \forall i \in \mathcal{I} \right\}. \quad (5)$$

We will slightly abuse notation by using symbol \mathbf{x} for a generic element of $\mathbb{R}_+^{|\mathcal{K}|}$; while $\mathbf{x}^r(\infty)$, $\mathbf{x}(\infty)$, $\mathbf{x}^r(t)$, and later $\mathbf{x}(t)$, refer to random vectors taking values in $\mathbb{R}_+^{|\mathcal{K}|}$.

Note that the asymptotic regime and the associated basic properties (3) and (4) hold *for any placement algorithm*. Indeed, (3) and (4) only depend on the fact that $Y_i^r(\infty)$ are mutually independent Poisson random variables with means $\rho_i r$, discussed earlier.

Consider the following linear program (LP) of minimizing $\sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}}$, the number of occupied servers on the fluid scale.

$$\text{(LP)} \quad \text{Minimize} \quad \sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}} \quad (6)$$

$$\text{subject to} \quad \sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}} = \rho_i \quad \forall i \in \mathcal{I}, \quad (7)$$

$$x_{\mathbf{k}} \geq 0, \quad \forall \mathbf{k} \in \mathcal{K}. \quad (8)$$

Denote by $\mathcal{X}^* \subseteq \mathcal{X}$ the set of its optimal solutions, and by L^* its optimal value. Since constraints (7) hold for $\mathbf{x}(\infty)$ under any placement algorithm, the optimal value L^* provides a lower bound on $\sum_{\mathbf{k}} x_{\mathbf{k}}(\infty)$, the steady-state total number of occupied servers on the fluid scale, under any placement algorithm.

3 Main Result

The **main result** of this paper is the following theorem. It shows that $\text{GRAND}(Z^p)$ is asymptotically optimal, when parameter p is sufficiently close to 1, which depends on the structure of the packing constraint specified by \mathcal{K} .

Theorem 4. Denote $\kappa \doteq 1 + \max_{\mathbf{k}} \sum_i k_i$, and assume that $p < 1$ and

$$1 - \kappa(1 - p) > 7/8, \quad \text{or equivalently, } p > 1 - 1/(8\kappa). \quad (9)$$

Consider a sequence of systems, with parameter $r \rightarrow \infty$, operating under the $\text{GRAND}(Z^p)$ algorithm. For each r , let $\mathbf{x}^r(\infty)$ denote the random state of the fluid-scaled process in the stationary regime. Then as $r \rightarrow \infty$,

$$d(\mathbf{x}^r(\infty), \mathcal{X}^*) \Rightarrow 0.$$

We would like to compare our main result, Theorem 4, with the main results (Theorems 3 and 4) of [23], on the asymptotic performance of the $\text{GRAND}(aZ)$ algorithm. For any given p that satisfies (9), $\text{GRAND}(Z^p)$ is asymptotically optimal in the sense that the *fluid-scaled optimality gap*, as measured by $d(\mathbf{x}^r(\infty), \mathcal{X}^*)$, goes to 0 as the system scale r grows to infinity. In contrast, $\text{GRAND}(aZ)$ is asymptotically optimal in the following (essentially, weaker) sense: for any fixed parameter $a > 0$ and under $\text{GRAND}(aZ)$, as $r \rightarrow \infty$, the fluid-scaled optimality gap $d(\mathbf{x}^r(\infty), \mathcal{X}^*)$ converges to a constant $c(a)$, which is not necessarily (and typically is not) zero; however, $c(a) \rightarrow 0$ as $a \rightarrow 0$. So, speaking informally, $\text{GRAND}(Z^p)$ optimality requires only that $r \rightarrow \infty$, while the optimality of $\text{GRAND}(aZ)$ requires that $r \rightarrow \infty$ and then $a \rightarrow 0$.

4 Basic System Dynamics

We now describe the basic system dynamics under customer arrivals and departures, and introduce some additional notation and terminology that will be used in the sequel. We note that the dynamics to be described here is general, and is not limited to $\text{GRAND}(Z^p)$, or any other placement algorithm. In other words, here the question of *how* placement decisions are made is irrelevant.

Define $\mathcal{M} = \{(\mathbf{k}, i) : \mathbf{k} \in \mathcal{K}, \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}\}$. Note that $(\mathbf{k}, i) \in \mathcal{M}$ are in one-to-one correspondence with the pairs $(\mathbf{k}, \mathbf{k} - \mathbf{e}_i)$ with $\mathbf{k} \in \mathcal{K}$ and $\mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}$, so (\mathbf{k}, i) can be viewed as a shorthand for $(\mathbf{k}, \mathbf{k} - \mathbf{e}_i)$. For this reason, we call a pair $(\mathbf{k}, i) \in \mathcal{M}$ an *edge*.

The placement of a type- i arrival into a server with configuration $\mathbf{k} - \mathbf{e}_i$ to form configuration \mathbf{k} is called an *arrival along the edge* (\mathbf{k}, i) , and the departure of a type- i customer from a server with configuration \mathbf{k} , which changes the server configuration to $\mathbf{k} - \mathbf{e}_i$, is called a *departure along the edge* (\mathbf{k}, i) .

If each non-empty server is viewed as a *particle*, then $X_{\mathbf{k}}$ is the number of particles with configuration \mathbf{k} . Then, arrivals and departures along edges can be interpreted as *particle movements* as follows. An arrival along an edge $(\mathbf{k}, i) \in \mathcal{M}$ at time t causes the following changes to the system state $\mathbf{X}(t)$. If $\mathbf{k} - \mathbf{e}_i = \mathbf{0}$, this means that the customer is placed into an empty server, and then the only state change is a *particle creation* at $X_{\mathbf{e}_i}$; i.e., $X_{\mathbf{e}_i}(t+) = X_{\mathbf{e}_i}(t-) + 1$, and $X_{\mathbf{k}'}(t+) = X_{\mathbf{k}'}(t-)$ for all $\mathbf{k}' \neq \mathbf{e}_i, \mathbf{k}' \in \mathcal{K}$. If $\mathbf{k} - \mathbf{e}_i \neq \mathbf{0}$, then the only state change is a *particle movement* from $X_{\mathbf{k} - \mathbf{e}_i}$ to $X_{\mathbf{k}}$; i.e., $X_{\mathbf{k} - \mathbf{e}_i}(t+) = X_{\mathbf{k} - \mathbf{e}_i}(t-) - 1$, $X_{\mathbf{k}}(t+) = X_{\mathbf{k}}(t-) + 1$, and $X_{\mathbf{k}'}(t+) = X_{\mathbf{k}'}(t-)$ for all $\mathbf{k}' \in \mathcal{K}$ such that $\mathbf{k}' \neq \mathbf{k}$ and $\mathbf{k}' \neq \mathbf{k} + \mathbf{e}_i$. Similarly, for a departure along an edge $(\mathbf{k}, i) \in \mathcal{M}$ at time t , if $\mathbf{k} = \mathbf{e}_i$, then we have a *particle annihilation* at $X_{\mathbf{k}}$; i.e., $X_{\mathbf{k}}(t+) = X_{\mathbf{k}}(t-) - 1$, and $X_{\mathbf{k}'}(t+) = X_{\mathbf{k}'}(t-)$ for all $\mathbf{k}' \neq \mathbf{k}$. If $\mathbf{k} \neq \mathbf{e}_i$, then we have a *particle movement* from $X_{\mathbf{k}}$ to $X_{\mathbf{k} - \mathbf{e}_i}$; i.e., $X_{\mathbf{k}}(t+) = X_{\mathbf{k}}(t-) - 1$, $X_{\mathbf{k} - \mathbf{e}_i}(t+) = X_{\mathbf{k} - \mathbf{e}_i}(t-) + 1$, and $X_{\mathbf{k}'}(t+) = X_{\mathbf{k}'}(t-)$ for all $\mathbf{k}' \in \mathcal{K}$ such that $\mathbf{k}' \neq \mathbf{k}$ and $\mathbf{k}' \neq \mathbf{k} - \mathbf{e}_i$.

5 Preliminaries

In this section, we present some definitions and facts from [23], to provide enough background needed for the proof of Theorem 4, our main result.

Dual characterization of (LP). Using the monotonicity of $\bar{\mathcal{K}}$ (cf. Definition 1), it is easy to check that if, in the program (LP) defined by (6)-(8), we replace equality constraints (7) with the inequality constraints

$$\sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}} \geq \rho_i, \quad \forall i, \quad (10)$$

we form a new linear program (LP') with the same optimal value as (LP). More explicitly, (LP') is given by

$$\begin{aligned} \text{(LP')} \quad & \text{Minimize } \sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}} \\ & \text{subject to } \sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}} \geq \rho_i \quad \forall i \in \mathcal{I}, \\ & \quad \quad \quad x_{\mathbf{k}} \geq 0 \quad \forall \mathbf{k} \in \mathcal{K}. \end{aligned}$$

Let \mathcal{X}^{**} denote the set of optimal solutions of (LP'). Then, \mathcal{X}^{**} contains \mathcal{X}^* , the set of optimal solutions of (LP); or more precisely, $\mathcal{X}^* = \mathcal{X}^{**} \cap \mathcal{X}$. The dual program (DUAL') of (LP') is given by

$$\text{(DUAL')} \quad \text{Maximize } \sum_{i \in \mathcal{I}} \rho_i \eta_i \quad (11)$$

$$\text{subject to } \sum_{i \in \mathcal{I}} k_i \eta_i \leq 1, \quad \forall \mathbf{k} \in \mathcal{K}, \quad (12)$$

$$\eta_i \geq 0, \quad \forall i \in \mathcal{I}. \quad (13)$$

The following lemma is a simple consequence of the Kuhn-Tucker theorem.

Lemma 5. *Vector \mathbf{x} is an optimal solution of (LP), i.e., $\mathbf{x} \in \mathcal{X}^*$, if and only if $\mathbf{x} \in \mathcal{X}$, and there exists a vector $\boldsymbol{\eta} = \{\eta_i, i \in \mathcal{I}\}$ that satisfies (12), (13) and the complementary slackness condition:*

$$\text{for any } \mathbf{k} \in \mathcal{K}, \text{ if } \sum_i k_i \eta_i < 1, \text{ then } x_{\mathbf{k}} = 0. \quad (14)$$

(Clearly, any such vector $\boldsymbol{\eta}$ is an optimal solution of (DUAL').)

Let \mathcal{H}^* be the set of vectors $\boldsymbol{\eta}$ that satisfy (12)-(14) for some $\mathbf{x} \in \mathcal{X}$.

A useful Lyapunov function. For any parameter $a \in (0, 1)$, define the function $L^{(a)} : \mathbb{R}_+^{|\mathcal{K}|} \rightarrow \mathbb{R}$ by

$$L^{(a)}(\mathbf{x}) \doteq -\frac{1}{\log a} \sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}} \log \left(\frac{x_{\mathbf{k}} c_{\mathbf{k}}}{ea} \right), \quad (15)$$

where $c_{\mathbf{k}} \doteq \prod_i k_i!$, $0! = 1$, and we use the convention that $0 \log 0 = 0$.

Function $L^{(a)}(\mathbf{x})$ can be viewed as an approximation to the linear function $\sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}}$. Indeed, it has been shown [23] that as $a \rightarrow 0$, $L^{(a)}(\mathbf{x}) \rightarrow \sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}}$ uniformly over any compact set.

Throughout the paper, we will use notation $b = -\log a$. Then, for each $\mathbf{k} \in \mathcal{K}$, we have

$$\frac{\partial}{\partial x_{\mathbf{k}}} L^{(a)}(\mathbf{x}) = \frac{1}{b} \log \left(\frac{c_{\mathbf{k}} x_{\mathbf{k}}}{a} \right). \quad (16)$$

Note that if we adopt the convention that

$$\frac{\partial}{\partial \mathbf{x}_0} L^{(a)}(\mathbf{x}) \Big|_{x_0=a} = 0, \quad (17)$$

then (16) is valid for $\mathbf{k} = \mathbf{0}$ and $x_0 = a$, which will be useful later.

The function $L^{(a)}$ is strictly convex in $\mathbf{x} \in \mathbb{R}_+^{|\mathcal{K}|}$. Consider the problem $\min_{\mathbf{x} \in \mathcal{X}} L^{(a)}(\mathbf{x})$. It is the following convex optimization problem (CVX(a)):

$$(CVX(a)) \quad \text{Minimize } L^{(a)}(\mathbf{x}) \quad (18)$$

$$\text{subject to } \sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}} = \rho_i, \quad \forall i \in \mathcal{I}, \quad (19)$$

$$x_{\mathbf{k}} \geq 0, \quad \forall \mathbf{k} \in \mathcal{K}. \quad (20)$$

Denote by $\mathbf{x}^{*,a} \in \mathcal{X}$ its unique optimal solution. The following lemma provides a crisp characterization of the point $\mathbf{x}^{*,a}$.

Lemma 6. *A point $\mathbf{x} \in \mathcal{X}$ is the optimal solution to (CVX(a)) defined by (18)-(20), i.e., $\mathbf{x} = \mathbf{x}^{*,a}$, if and only if it has a product form representation*

$$x_{\mathbf{k}}^{*,a} = \frac{a}{c_{\mathbf{k}}} \exp \left[b \sum_i k_i \nu_i^{*,a} \right] = \frac{1}{c_{\mathbf{k}}} a^{1 - \sum_i k_i \nu_i^{*,a}}, \quad \forall \mathbf{k} \in \mathcal{K}, \quad (21)$$

for some vector $\boldsymbol{\nu}^{*,a} = \{\nu_i^{*,a}, i \in \mathcal{I}\}$.

Proof sketch. For each $i \in \mathcal{I}$, let ν_i be the dual variable that corresponds to the equality constraint $\sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}} = \rho_i$, and consider the Lagrangian

$$L^{(a)}(\mathbf{x}) + \sum_i \nu_i \left(\rho_i - \sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}} \right).$$

By setting partial derivatives of the Lagrangian to zero, we get

$$\frac{1}{b} \log \left[\frac{x_{\mathbf{k}} c_{\mathbf{k}}}{a} \right] - \sum_i \nu_i k_i = 0, \quad \forall \mathbf{k} \in \mathcal{K}.$$

$\mathbf{x}^{*,a} \in \mathcal{X}$ is the optimal solution to (CVX(a)) if and only if there exist Lagrange multipliers $\boldsymbol{\nu}^{*,a} = \{\nu_i^{*,a}, i \in \mathcal{I}\}$ such that

$$\frac{1}{b} \log \left[\frac{x_{\mathbf{k}}^{*,a} c_{\mathbf{k}}}{a} \right] - \sum_i \nu_i^{*,a} k_i = 0, \quad \forall \mathbf{k} \in \mathcal{K},$$

which leads to the product form representation (21). □

A simple consequence of Lemma 6 is that the Lagrange multipliers $\nu_i^{*,a}$ are unique and are equal to $1 - \log(x_{\mathbf{e}_i}^{*,a}) / \log a$, by considering (21) for \mathbf{e}_i , $i \in \mathcal{I}$. The following result is from [23].

Theorem 7. *Let $\mathbf{x}^{*,a}$ and $\boldsymbol{\nu}^{*,a}$ be as in Lemma 6. Then, as $a \downarrow 0$, $\mathbf{x}^{*,a} \rightarrow \mathcal{X}^*$ and $\boldsymbol{\nu}^{*,a} \rightarrow \mathcal{H}^*$.*

Note that the convergence $\mathbf{x}^{*,a} \rightarrow \mathcal{X}^*$ is an easy consequence of the fact that $L^{(a)}(\mathbf{x}) - \sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}} \rightarrow 0$ u.o.c., as $a \rightarrow 0$.

In [23], we considered *fluid limit* trajectories $\mathbf{x}(\cdot)$, which, speaking informally, arise as limits of the fluid-scaled process $\mathbf{x}^r(\cdot)$ as $r \rightarrow \infty$. We showed there that for any fixed $a > 0$, under GRAND(aZ), the derivative $(d/dt)L^{(a)}(\mathbf{x}(t))$ along a fluid limit trajectory $\mathbf{x}(\cdot)$ is always negative, whenever $\mathbf{x}(t) \neq \mathbf{x}^{*,a}$; consequently,

$L^{(a)}$ served as a Lyapunov function, which allowed us to prove that $\mathbf{x}(t) \rightarrow \mathbf{x}^{*,a}$, as $t \rightarrow \infty$, along any fluid limit trajectory. This in turn was the key property that led to establishing the fact that the random steady-state system state $\mathbf{x}^r(\infty)$ concentrates on $\mathbf{x}^{*,a}$ as $r \rightarrow \infty$; together with Theorem 7 this means the asymptotic optimality of GRAND(aZ), if the limit $r \rightarrow \infty$ of steady states $\mathbf{x}^r(\infty)$ is taken first, and the limit $a \rightarrow 0$ is taken after that. In this paper we will also make use of the family of Lyapunov functions $L^{(a)}(\mathbf{x})$ defined in (15), but use them in a different way, because the analysis of GRAND(Z^p) cannot (as will be explained later) rely on an analysis of fluid-limit trajectories.

We continue to introduce some expressions and their properties, derived in [23] and which are closely related to the Lyapunov function $L^{(a)}$. In [23], the expression $\Xi(\mathbf{x})$ defined below in (25) is the derivative of the Lyapunov function $L^{(a)}(\mathbf{x}(t))$ with respect to time, along a fluid limit trajectory $\mathbf{x}(\cdot)$. In the context of this paper, Ξ will be interpreted and used differently; roughly speaking, it will be the value of the process generator when applied to $L^{(a)}$.

Consider $\mathbf{x} \in \mathbb{R}_+^{|\mathcal{K}|}$ such that $x_{\mathbf{k}} > 0$, $\forall \mathbf{k} \in \mathcal{K}$. For an ordered pair of edges (\mathbf{k}, i) and (\mathbf{k}', i) we define

$$\xi_{\mathbf{k}, \mathbf{k}', i} = \frac{1}{b} \left[\log(k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'}) - \log(k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i}) \right] \frac{k_i \mu_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i}}{x_{(i)}},$$

where *by convention*,

$$\mathbf{x}_0 = a \quad \text{and, correspondingly,} \quad x_{(i)} \doteq \mathbf{x}_0 + \sum_{\mathbf{k} \in \mathcal{K}: \mathbf{k}+e_i \in \mathcal{K}} x_{\mathbf{k}}, \quad (22)$$

and $\xi_{\mathbf{k}, \mathbf{k}', i}$ is defined for cases when either $\mathbf{k} - e_i = \mathbf{0}$ or $\mathbf{k}' - e_i = \mathbf{0}$ as well. Let us note that in the expression of $\xi_{\mathbf{k}, \mathbf{k}', i}$,

$$\begin{aligned} & \frac{1}{b} \left[\log(k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'}) - \log(k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i}) \right] = \\ & \left[\frac{\partial}{\partial x_{\mathbf{k}'}} L^{(a)}(\mathbf{x}) - \frac{\partial}{\partial x_{\mathbf{k}'-e_i}} L^{(a)}(\mathbf{x}) \right] - \left[\frac{\partial}{\partial x_{\mathbf{k}}} L^{(a)}(\mathbf{x}) - \frac{\partial}{\partial x_{\mathbf{k}-e_i}} L^{(a)}(\mathbf{x}) \right]. \end{aligned} \quad (23)$$

We have

$$\xi_{\mathbf{k}, \mathbf{k}', i} + \xi_{\mathbf{k}', \mathbf{k}, i} = \frac{\mu_i}{b x_{(i)}} \left[\log(k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'}) - \log(k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i}) \right] [k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i} - k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'}] \leq 0, \quad (24)$$

and the inequality is strict unless $k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'} = k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i}$. Finally, we define

$$\Xi(\mathbf{x}) = \sum_i \sum_{\mathbf{k}, \mathbf{k}'} [\xi_{\mathbf{k}, \mathbf{k}', i} + \xi_{\mathbf{k}', \mathbf{k}, i}]. \quad (25)$$

Observe that $\Xi(\mathbf{x}) < 0$ unless \mathbf{x} has a product form representation (21), for some vector $\boldsymbol{\nu}^{*,a}$. This is because $\Xi(\mathbf{x}) = 0$ if and only if the equality in (24) holds for all pairs of edges $\{(\mathbf{k}, i), (\mathbf{k}', i)\}$, which is true if and only if $k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'} = k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i}$ for all pairs $\{(\mathbf{k}, i), (\mathbf{k}', i)\}$, a condition that is equivalent to a product form presentation (21) of \mathbf{x} , for some $\boldsymbol{\nu}^{*,a}$. Let us also note that here the vector $\boldsymbol{\nu}^{*,a}$ need not correspond to the (unique) Lagrange multipliers of the program (CVX(a)), because \mathbf{x} does not necessarily satisfy the equalities (19); i.e., we do not necessarily have $\mathbf{x} \in \mathcal{X}$ (recall the definition of \mathcal{X} in (5)). If $\mathbf{x} \in \mathcal{X}$, we will have $\mathbf{x} = \mathbf{x}^{*,a}$, exactly as in (21).

We now provide a more concrete interpretation of (25). Consider a state \mathbf{x} and suppose that the time derivatives (or rates of changes) of \mathbf{x} are defined by the following dynamics. Along each edge (\mathbf{k}, i) , “mass” is moving from $x_{\mathbf{k}}$ to $x_{\mathbf{k}-e_i}$ at the rate

$$\tilde{w}_{\mathbf{k}, i} = k_i \mu_i x_{\mathbf{k}}. \quad (26)$$

(Value $\tilde{w}_{\mathbf{k}, i}$ can be interpreted as the *fluid-scaled* rate of type- i departures along the edge (\mathbf{k}, i) .) In addition, along each edge (\mathbf{k}, i) , “mass” is moving from $x_{\mathbf{k}-e_i}$ to $x_{\mathbf{k}}$ at the rate

$$\tilde{v}_{\mathbf{k}, i} = \frac{\tilde{\lambda}_i x_{\mathbf{k}-e_i}}{x_{(i)}}, \quad (27)$$

where we denote

$$\tilde{\lambda}_i = \sum_{\mathbf{k} \in \mathcal{K}} k_i \mu_i x_{\mathbf{k}}, \quad (28)$$

and $x_{(i)}$ is defined by (22). (Values $\tilde{v}_{\mathbf{k},i}$ and $\tilde{\lambda}_i$ defined in (27) and (28) can be interpreted as follows. If $\mathbf{x} \in \mathcal{X}$, then $\tilde{\lambda}_i = \lambda_i$ for all $i \in \mathcal{I}$, by the equality constraints (19), and $\tilde{v}_{\mathbf{k},i}$ of (27) is the *fluid-scaled* rate of type- i arrivals along the edge (\mathbf{k}, i) , under the GRAND(Z^p) algorithm. In general, this is not necessarily the case, since we consider a generic system state \mathbf{x} , which need not belong to \mathcal{X} .)

Equations (26) and (27) define the derivatives of $x_{\mathbf{k}}$, for all $\mathbf{k} \in \mathcal{K}$. By convention, $x_{\mathbf{0}}$ remains constant at a , so it has zero derivative. With such derivatives of \mathbf{x} and $x_{\mathbf{0}}$, given the convention (22) and recalling (23), it is not difficult to check that $\Xi(\mathbf{x})$ is the time derivative of $L^{(a)}(\mathbf{x})$:

$$\Xi(\mathbf{x}) = \sum_{(\mathbf{k},i) \in \mathcal{M}} \left[\frac{\partial}{\partial x_{\mathbf{k}}} L^{(a)}(\mathbf{x}) - \frac{\partial}{\partial x_{\mathbf{k}-e_i}} L^{(a)}(\mathbf{x}) \right] [\tilde{v}_{\mathbf{k},i} - \tilde{w}_{\mathbf{k},i}], \quad (29)$$

where the convention (17) is used.

For future reference, we will use the following notation for the term $\log(k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'}) - \log(k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i})$ in (24), defined for a pair of edges (\mathbf{k}, i) and (\mathbf{k}', i) :

$$\chi_{\mathbf{k},\mathbf{k}',i}(\mathbf{x}) = \log(k'_i x_{\mathbf{k}-e_i} x_{\mathbf{k}'}) - \log(k_i x_{\mathbf{k}} x_{\mathbf{k}'-e_i}). \quad (30)$$

6 Proof of Theorem 4

This section is organized as follows. In Subsection 6.1, we present the Lyapunov function that we will use for each system indexed by r , and derive some basic properties of this Lyapunov function. In particular, we will see that the Lyapunov function can be thought of as an approximation to the linear objective, $\sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}}$. We also point out the main difficulties in establishing the asymptotic optimality of GRAND(Z^p) here. In Subsection 6.2, we present some high probability bounds of the process $\mathbf{x}(\cdot)$ in the stationary regime, for large r . Making use of the high-probability bounds in Subsection 6.2, we derive high-probability estimates on the change of the Lyapunov function in Subsection 6.3, and use these estimates to conclude the proof of Theorem 4 in Subsection 6.4.

Throughout this section, we drop superscript r in notation – such as that appears in $\mathbf{x}^r(t)$ – pertaining to processes with parameter r .

6.1 A Lyapunov Function and Some Basic Properties

We will use the Lyapunov function $L^{(a)}$ for the proof of Theorem 4, where the parameter a depends on the system scale r as $a = r^{p-1}$. Correspondingly, b will also depend on r , with $b = -\log a = (1-p)\log r$. When it is clear from the context that we are working with $L^{(a)}$ for a given fixed r (and corresponding $a = r^{p-1}$), we often drop the superscript (a) .

Let us now provide an important remark regarding the Lyapunov function $L^{(a)}$ with $a = r^{p-1}$, and the GRAND(Z^p) algorithm. On the one hand, for any given system scale parameter r , $L^{(a)}$ has a fixed parameter $a = r^{p-1}$. Furthermore, the expressions (29) and (25) for $\Xi(\mathbf{x})$, the time derivative of $L^{(a)}$, make use of the *convention* (22), where $a = r^{p-1}$. On the other hand, GRAND(Z^p) itself does *not* use or need to know the parameter r , since the actual number of zero-servers that it uses at time t is $X_{\mathbf{0}}(t) = Z^p(t)$, which only uses the knowledge of $Z(t)$, the total number of customers at time t .

At this point, we can highlight the key technical difficulty of the analysis of GRAND(Z^p), compared to that of GRAND(aZ) in [23]. Under GRAND(aZ), the fluid-scaled quantities $x_{\mathbf{k}}^r$ are $O(1)$ for *all* configurations

\mathbf{k} . Consequently, as $r \rightarrow \infty$, the fluid-scaled process $\mathbf{x}^r(\cdot)$ converges to a well-defined fluid limit $\mathbf{x}(\cdot)$, for any well-defined limiting initial state $\mathbf{x}(0)$. The Lyapunov function $L^{(a)}$, with fixed $a > 0$, is then used in [23] to show the convergence of all fluid trajectories $\mathbf{x}(t)$ to the point $\mathbf{x}^{*,a}$, as $t \rightarrow \infty$, which is the key part of the analysis of $\text{GRAND}(aZ)$. Therefore, in [23], it sufficed to work with fluid limits and derivatives of the Lyapunov function along fluid limit trajectories.

In this paper, to analyze $\text{GRAND}(Z^p)$, we use $L^{(r^{p-1})}$ as a Lyapunov function for a system with given r . The key difficulty is that under $\text{GRAND}(Z^p)$, some of the fluid-scaled $x_{\mathbf{k}}^r$ are $o(1)$; for example, $x_{\mathbf{0}}^r = O(r^{p-1})$. A further complication is that different $x_{\mathbf{k}}^r$ are on *different scales* with respect to r ; namely, they are $O(r^{s(\mathbf{k})})$ with the power $s(\mathbf{k})$ depending on \mathbf{k} , $s(\mathbf{k}) \in (0, 1]$. Because of the multi-scale system dynamics, the conventional fluid limit is not sufficient to establish a negative drift of the Lyapunov function. Moreover, the use of other, *local*, fluid limits (obtained under different space/time scalings) also does not appear to be particularly useful, because, again, $x_{\mathbf{k}}^r$ evolve on different scales. To overcome this difficulty, in Section 6.3, we use direct (and more involved) probabilistic estimates of the Lyapunov function increments, which are based on martingale theory, and which do not involve fluid or local fluid limits.

The following lemmas state some elementary properties related to the Lyapunov function $L^{(a)}$ with $a = r^{p-1}$.

Lemma 8. *Let $\varepsilon \in (0, \frac{1}{2})$. Consider a sequence of points $\mathbf{x}^{\circ,a}$ indexed by $a = r^{p-1}$, which satisfy*

$$\left| \sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}}^{\circ,a} - \rho_i \right| \leq \frac{r^{-1/2+\varepsilon}}{I}, \quad \forall i, \text{ for sufficiently large } r, \quad (31)$$

and

$$\chi_{\mathbf{k}, \mathbf{k}', i}(\mathbf{x}^{\circ,a}) \rightarrow 0 \text{ for all pairs of edges } (\mathbf{k}, i), (\mathbf{k}', i),$$

where we recall the definition of $\chi_{\mathbf{k}, \mathbf{k}', i}$ in (30). Then, as $r \rightarrow \infty$, $L^{(a)}(\mathbf{x}^{\circ,a}) \rightarrow L^*$, where L^* is the optimal value of (LP) defined by (6)–(8).

Proof. For each r , and the corresponding $a = r^{p-1}$, denote

$$\nu_i^{*,a} = 1 - \frac{\log(x_{\mathbf{e}_i}^{\circ,a})}{\log a}, \quad \forall i \in \mathcal{I}.$$

Define $\mathbf{x}^{*,a} = \{x_{\mathbf{k}}^{*,a}, \mathbf{k} \in \mathcal{K}\}$ via product form (21) corresponding to these $\nu_i^{*,a}$:

$$x_{\mathbf{k}}^{*,a} = \frac{a}{c_{\mathbf{k}}} \exp \left[b \sum_i k_i \nu_i^{*,a} \right] = \frac{1}{c_{\mathbf{k}}} a^{1 - \sum_i k_i \nu_i^{*,a}}, \quad \mathbf{k} \in \mathcal{K}. \quad (32)$$

Note that for any r (and correspondingly, $a = r^{p-1}$), $\mathbf{x}^{*,a}$ and $\boldsymbol{\nu}^{*,a}$ need not be the optimal primal and dual solutions to $(\text{CVX}(a))$, defined by (18)–(20), since $\mathbf{x}^{*,a}$ need not belong to \mathcal{X} . Instead, point $\mathbf{x}^{*,a}$ is the optimal solution for the problem

$$\text{minimize } L^{(a)}(\mathbf{x}) \quad (33)$$

$$\text{subject to } \sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}} = \sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}}^{*,a}, \quad \forall i \in \mathcal{I}, \quad (34)$$

$$x_{\mathbf{k}} \geq 0, \quad \forall \mathbf{k} \in \mathcal{K}. \quad (35)$$

We now claim that for all $\mathbf{k} \in \mathcal{K}$,

$$x_{\mathbf{k}}^{\circ,a} / x_{\mathbf{k}}^{*,a} \rightarrow 1 \text{ as } r \rightarrow \infty. \quad (36)$$

By convention, $x_{\mathbf{0}}^{\circ,a} = x_{\mathbf{0}}^{*,a} = a = r^{p-1}$. It is also easy to check that for all $i \in \mathcal{I}$, $x_{\mathbf{e}_i}^{\circ,a} = x_{\mathbf{e}_i}^{*,a}$. Consider a general configuration $\mathbf{k} \in \mathcal{K}$. For an edge (\mathbf{k}, i) , we use $\tilde{\chi}_{(\mathbf{k}, i)}$ as a shorthand for $\chi_{\mathbf{e}_i, \mathbf{k}, i}(\mathbf{x}^{\circ,a})$, and we have

$$\begin{aligned} \tilde{\chi}_{(\mathbf{k}, i)} &= \chi_{\mathbf{e}_i, \mathbf{k}, i}(\mathbf{x}^{\circ,a}) \\ &= \log(k_i x_{\mathbf{0}}^{\circ,a} x_{\mathbf{k}}^{\circ,a}) - \log(x_{\mathbf{e}_i}^{\circ,a} x_{\mathbf{k} - \mathbf{e}_i}^{\circ,a}) \\ &= \log(k_i a x_{\mathbf{k}}^{\circ,a}) - \log(x_{\mathbf{e}_i}^{*,a} x_{\mathbf{k} - \mathbf{e}_i}^{\circ,a}). \end{aligned}$$

By some simple algebraic manipulation, we have the recursion

$$x_{\mathbf{k}}^{\circ,a} = \frac{e^{\tilde{\chi}(\mathbf{k},i)} x_{\mathbf{e}_i}^{*,a} x_{\mathbf{k}-\mathbf{e}_i}^{\circ,a}}{k_i a} = \frac{1}{k_i} e^{\tilde{\chi}(\mathbf{k},i)} a^{-\nu_i^{*,a}} x_{\mathbf{k}-\mathbf{e}_i}^{\circ,a}. \quad (37)$$

Consider a path that connects configuration \mathbf{k} and $\mathbf{0}$ using the following sequence of edges: the first k_1 edges are $(\mathbf{k}, 1), (\mathbf{k} - \mathbf{e}_1, 1), \dots, (\mathbf{k} - (k_1 - 1)\mathbf{e}_1, 1)$, followed by k_2 edges of the form $(\mathbf{k} - (k_1 - 1)\mathbf{e}_1, 2), (\mathbf{k} - (k_1 - 1)\mathbf{e}_1 - \mathbf{e}_2, 2), \dots, (\mathbf{k} - (k_1 - 1)\mathbf{e}_1 - (k_2 - 1)\mathbf{e}_2, 2)$, etc. Thus, the path uses $\sum_i k_i$ edges to connect \mathbf{k} with $\mathbf{0}$. Using \tilde{e} to denote a generic edge that belongs to this path, and using the recursion (37), it is easy to see that

$$x_{\mathbf{k}}^{\circ,a} = \frac{1}{c_{\mathbf{k}}} \left[\prod_{\tilde{e}} \exp(\tilde{\chi}_{\tilde{e}}) \right] a^{1 - \sum_i k_i \nu_i^{*,a}} = \left[\prod_{\tilde{e}} \exp(\tilde{\chi}_{\tilde{e}}) \right] x_{\mathbf{k}}^{*,a}.$$

Thus,

$$\frac{x_{\mathbf{k}}^{\circ,a}}{x_{\mathbf{k}}^{*,a}} = \left[\prod_{\tilde{e}} \exp(\tilde{\chi}_{\tilde{e}}) \right] \rightarrow 1,$$

since $\tilde{\chi}_{\tilde{e}} \rightarrow 0$ for all edges \tilde{e} . This completes the proof of the claim.

By (36), we have

$$\|\mathbf{x}^{\circ,a} - \mathbf{x}^{*,a}\| \rightarrow 0, \quad r \rightarrow \infty.$$

In addition, using condition (31), we have that for each $i \in \mathcal{I}$, the term $\sum_{\mathbf{k} \in \mathcal{K}} k_i x_{\mathbf{k}}^{*,a}$ on the RHS of the equality constraint (34) converges to ρ_i as $r \rightarrow \infty$. Furthermore, recall that $|L^{(a)}(\mathbf{x}) - \sum_{\mathbf{k} \in \mathcal{K}} x_{\mathbf{k}}| \rightarrow 0$, uniformly on \mathbf{x} from a bounded set. Therefore, we must have $L^{(a)}(\mathbf{x}^{*,a}) \rightarrow L^*$. \square

Lemma 8 implies the following property that we will actually use.

Lemma 9. *For any $\gamma > 0$, there exists $\delta_1 > 0$ such that, for all sufficiently large r , condition (31) and the condition that*

$$|\chi_{\mathbf{k},\mathbf{k}',i}| \leq \delta_1 \quad \text{for all pairs of edges } (\mathbf{k}, i), (\mathbf{k}', i)$$

imply

$$L^{(a)}(\mathbf{x}) - L^* \leq \gamma.$$

6.2 High Probability Steady-State Bounds on $\mathbf{x}(\cdot)$

From now on in the paper, we fix p satisfying (9), and a corresponding s , satisfying $7/8 < s < 1 - \kappa(1 - p)$.

We often refer to the following conditions, at a given time $t \geq 0$, with some constants $c > 0$ and $\varepsilon > 0$:

$$|Z(t)/r - 1| \leq r^{-1/2+\varepsilon}; \quad (38)$$

$$\left| \sum_{\mathbf{k}} k_i x_{\mathbf{k}}(t) - \rho_i \right| \leq \frac{r^{-1/2+\varepsilon}}{I}, \quad \forall i; \quad (39)$$

$$x_{\mathbf{k}}(t) \geq cr^{s-1}, \quad \forall \mathbf{k}. \quad (40)$$

It is easy to see that (39) implies (38).

Lemma 10. *Let $\alpha > 0$ and $\varepsilon > 0$. Consider our system in the stationary regime for each r . Then,*

$$\mathbb{P}\{\text{condition (39) holds for all } t \in [0, r^\alpha]\} \rightarrow 1, \quad r \rightarrow \infty.$$

The proof of Lemma 10 is provided in Appendix A. The statement of Lemma 10 is very intuitive, because we know that in the stationary regime, $Y_i(t)$ has Poisson distribution with mean $\rho_i r$ for any t . Since $Y_i(t) = \sum_{\mathbf{k} \in \mathcal{K}} k_i X_{\mathbf{k}}(t)$, condition (39) clearly holds for any given t . The proof shows that, in fact, with high probability, (39) holds on any r^α -long time interval.

Lemma 11. *Let $\alpha > 0$. Consider our system in the stationary regime for each r . Then, there exists $c > 0$ such that*

$$\mathbb{P}\{\text{condition (40) holds for } t \in [0, r^\alpha]\} \rightarrow 1, \quad r \rightarrow \infty.$$

The proof of Lemma 11 is provided in Appendix B, where we establish a stronger property: with high probability, there exists some $c > 0$ such that for all $\mathbf{k} \in \mathcal{K}$ and for all $t \in [0, r^\alpha]$, $x_{\mathbf{k}}(t) \geq cr^{s(\mathbf{k})-1}$, where $s(\mathbf{k}) = 1 - (\sum_i k_i + 1)(1-p) \geq 1 - \kappa(1-p) > s$. Here we provide a sketch of the argument used to establish this stronger property. By Lemma 10, with high probability, the fluid-scaled total number of customers, $Z(t)/r$, is in steady state close to 1 for all $t \in [0, r^\alpha]$. Therefore, $x_{\mathbf{0}}$ is close to $Z^p/r \approx r^{p-1}$, and the property should hold for $\mathbf{k} = \mathbf{0}$. Then, we use an induction argument. Fix configuration $\mathbf{k} \neq \mathbf{0}$, and suppose the property is true for all configurations $\mathbf{k}' \neq \mathbf{k}$, $\mathbf{k}' \leq \mathbf{k}$. We argue that it should then hold for \mathbf{k} and some c (possibly rechosen to be smaller), as follows. Pick $\mathbf{k}' = \mathbf{k} - \mathbf{e}_i$ for some i ; such \mathbf{k}' exists since $\mathbf{k} \neq \mathbf{0}$. We account separately for transitions that increase $x_{\mathbf{k}}$ and those that decrease $x_{\mathbf{k}}$. First, type- i arrivals along edge (\mathbf{k}, i) increase $x_{\mathbf{k}}$, and take place at (unscaled) rate $\lambda_i r x_{\mathbf{k}'}/x_{(i)}$. Furthermore, this arrival process is lower bounded by a Poisson process of (unscaled) rate $c_1 r \left(r^{s(\mathbf{k}')}/r \right) = c_1 r^{s(\mathbf{k}')}$, for a constant $c_1 > 0$, because $X_{(i)}/r \leq (X_{\mathbf{0}} + Z)/r$ is essentially upper bounded by 1. Second, transitions that would decrease $x_{\mathbf{k}}$ consist of departures from \mathbf{k} and arrivals into \mathbf{k} . When $x_{\mathbf{k}}(t) \leq cr^{s(\mathbf{k})-1}$, the (unscaled) rates of these transitions are upper bounded by $c_2 cr^{s(\mathbf{k})}$ and $c_3 cr^{(r^{s(\mathbf{k})}/r)^p} = c_3 cr^{s(\mathbf{k})+1-p}$, respectively, for some constants $c_2 > 0$ and $c_3 > 0$. Thus, the transitions that would decrease $x_{\mathbf{k}}$ are upper bounded by a Poisson process with (unscaled) rate $c_2 cr^{s(\mathbf{k})} + c_3 cr^{s(\mathbf{k})+1-p} \leq c_4 cr^{s(\mathbf{k})+1-p} = c_4 cr^{s(\mathbf{k}')}$, for some $c_4 > 0$. Therefore, if we choose constant c sufficiently small, then when $x_{\mathbf{k}}(t) \leq cr^{s(\mathbf{k})}/r$, the rate of transitions that increase $x_{\mathbf{k}}$ dominates the rate of transitions that decrease it, at least by some factor greater than 1. Thus, with high probability, $x_{\mathbf{k}}(t)$ “should” stay above $cr^{s(\mathbf{k})-1}$.

6.3 Lyapunov Function Drift Estimates

In this subsection, we consider an artificial process for which (38)-(40) hold w.p.1 over a r^{s-1} -long interval, and establish high probability estimates of the change of the Lyapunov function $L^{(r^{p-1})}$, for the artificial process over that interval. In Subsection 6.4, we see that these estimates can be used for analyzing the change in $L^{(r^{p-1})}$ under the original process, since from Subsection 6.2, (38)-(40) hold with high probability, and the artificial and the original processes coincide with high probability. Let us proceed with details.

Consider a single r^{s-1} -long interval, and let $T = r^{s-1}$. The initial state is such that, at $t = 0$, conditions (38), (39), and (40) hold for some constant $c > 0$ and a small constant $\varepsilon > 0$, where the magnitude of ε will be specified later (at the end of this subsection).

Let $\widehat{\mathcal{X}} = \widehat{\mathcal{X}}^r$ denote the set of states \mathbf{x} that satisfy (38)-(40). For $\mathbf{x} \in \widehat{\mathcal{X}}$, we have:

$$\left| \frac{\partial}{\partial x_{\mathbf{k}}} L(\mathbf{x}) \right| \leq C_1, \quad (41)$$

$$0 \leq \left| \frac{\partial^2}{\partial x_{\mathbf{k}}^2} L(\mathbf{x}) \right| \leq C_2' [r^{s-1} \log r]^{-1} \leq C_2 r^{1-s}, \quad (42)$$

for some constants C_1, C_2 and all sufficiently large r .

Our process $\mathbf{x}(t)$ is a pure jump process – the jumps are caused by customer arrivals and departures. For the rest of this subsection, we consider an artificial version of the process, which we now describe. In this artificial process, for which, with some abuse, we keep the same notation $\mathbf{x}(\cdot)$, some of the jumps (i.e., transitions) are not allowed to occur. Specifically, if ξ is a point in time when an arrival or departure takes place, we call it a point of *potential transition*. If this potential transition keeps the state within conditions (38)-(40), we allow it to occur, so this becomes an actual transition for the artificial process. If the potential transition were to take us to a state violating (38)-(40), we do *not* allow this transition to actually occur, and keep the state unchanged.

We now introduce some notation related to the artificial process. For $\mathbf{x} \in \widehat{\mathcal{X}}$, let $\mathcal{T}(\mathbf{x}) = \mathcal{T}^r(\mathbf{x})$ denote the set of all potential transitions out of it. Let $\mathcal{T} = \mathcal{T}^r = \cup_{\mathbf{x} \in \widehat{\mathcal{X}}} \mathcal{T}(\mathbf{x})$ be the set of all possible potential transitions from the states $\mathbf{x} \in \widehat{\mathcal{X}}$. Note that for any r , the cardinality of $\widehat{\mathcal{X}}$ and of \mathcal{T} are both finite. For any $\alpha \in \mathcal{T}(\mathbf{x})$, let \mathbf{x}_α denote the new state after potential transition α from \mathbf{x} . Denote by $\widehat{\mathcal{T}}(\mathbf{x}) \subseteq \mathcal{T}(\mathbf{x})$ the set of all allowed potential transition out of \mathbf{x} , i.e., those for which $\mathbf{x}_\alpha \in \widehat{\mathcal{X}}$. Finally, let $\widehat{\mathcal{T}} = \cup_{\mathbf{x} \in \widehat{\mathcal{X}}} \widehat{\mathcal{T}}(\mathbf{x}) \subseteq \mathcal{T}$ denote the set of all potential transitions that are allowed, and for $\mathbf{x} \in \widehat{\mathcal{X}}$ and $\alpha \in \mathcal{T}(\mathbf{x})$, denote by β_α the (infinitesimal) rate of potential transition α from \mathbf{x} .

We can and do construct the artificial process on the following probability space. Let the initial state $\mathbf{x}(0) \in \widehat{\mathcal{X}}$ be fixed. For each potential transition type $\alpha \in \mathcal{T}$, there is an independent unit-rate Poisson process $\Pi_\alpha(\cdot)$. The counting process of potential α -transitions is

$$H_\alpha(t) = \Pi_\alpha \left(\int_0^t \beta_\alpha I\{\alpha \in \mathcal{T}(\mathbf{x}(\xi))\} d\xi \right).$$

Denote by $\tau_\alpha(t)$ the set of time points where jumps of $H_\alpha(\cdot)$ occurs in $[0, t]$; these are the potential transition time points in $[0, t]$. Denote by $\tau(t) = \cup_{\alpha \in \mathcal{T}} \tau_\alpha(t)$ and $\widehat{\tau}(t) = \cup_{\alpha \in \widehat{\mathcal{T}}} \tau_\alpha(t)$ the sets of all potential transition points and allowed (i.e., actual) transition points, respectively, in $[0, t]$. The state $\mathbf{x}(t)$ at time t is the state after the last allowed potential transition. More formally, if $\max \widehat{\tau}(t) \in \tau_\alpha(t)$ and $\alpha \in \mathcal{T}(\mathbf{x}')$, then $\mathbf{x}(t) = \mathbf{x}'_\alpha$; if $\widehat{\tau}(t)$ is empty, then $\mathbf{x}(t) = \mathbf{x}(0)$. It is easy to see that, w.p.1, this construction uniquely determines the realization of the process, including all potential transition points.

Consider the following pure-jump process:

$$F(t) = \sum_{\alpha \in \mathcal{T}} \sum_{s \in \tau_\alpha(t)} \Delta_\alpha \equiv \sum_{\alpha \in \mathcal{T}} \Delta_\alpha H_\alpha(t),$$

where

$$\Delta_\alpha = L(\mathbf{x}_\alpha) - L(\mathbf{x}), \quad \alpha \in \mathcal{T}(\mathbf{x}).$$

Note that the above definition of $F(t)$ has the sum over *all* potential transitions, both allowed and not allowed. We observe for future reference that, if the artificial process realization is such that all potential transitions in $[0, t]$ are allowed, i.e., they are actual transitions, then $F(t) = L(\mathbf{x}(t)) - L(\mathbf{x}(0))$.

Denote

$$AL(\mathbf{x}) = \sum_{\alpha \in \mathcal{T}(\mathbf{x})} \beta_\alpha \Delta_\alpha \equiv \sum_{\alpha \in \mathcal{T}(\mathbf{x})} \beta_\alpha [L(\mathbf{x}_\alpha) - L(\mathbf{x})].$$

Let us remark that for a process where all potential transitions are allowed, AL is the process generator, if L is within its domain. For our artificial process $\mathbf{x}(\cdot)$, however, $AL(\mathbf{x})$ is just a formally defined expression.

We have

$$F(t) = M(t) + B(t), \tag{43}$$

where

$$B(t) = \int_0^t AL(\mathbf{x}(s)) ds = \sum_{\alpha \in \mathcal{T}} B_\alpha(t), \tag{44}$$

$$B_\alpha(t) = \Delta_\alpha \bar{H}_\alpha(t), \quad \bar{H}_\alpha(t) \equiv \int_0^t \beta_\alpha I\{\alpha \in \mathcal{T}(\mathbf{x}(s))\} ds, \tag{45}$$

$$M(t) = \sum_{\alpha \in \mathcal{T}} M_\alpha(t), \tag{46}$$

$$M_\alpha(t) = \Delta_\alpha [\Pi_\alpha(\bar{H}_\alpha(t)) - \bar{H}_\alpha(t)]. \tag{47}$$

Note that $H_\alpha(t) = \Pi_\alpha(\bar{H}_\alpha(t))$. We have (from, e.g., Lemmas 3.1 and 3.2 in [19]), that each $M_\alpha(t)$, $t \geq 0$, is a martingale w.r.t. the filtration describing the history of the process up to time t . (This itself is a

martingale, there is no need to localize, as in [19], because in our case the total transition rate and all jump sizes are uniformly bounded.) Moreover, its predictable and optional quadratic variations are, respectively,

$$\langle M_\alpha \rangle(t) = \Delta_\alpha^2 \bar{H}_\alpha(t), \quad [M_\alpha](t) = \Delta_\alpha^2 H_\alpha(t).$$

We see that $M(\cdot)$ is also a martingale. Furthermore, since w.p.1 all potential transition time points are distinct, we have

$$[M](t) = \sum_{\alpha \in \mathcal{T}} [M_\alpha](t) = \sum_{\alpha \in \mathcal{T}} \Delta_\alpha^2 H_\alpha(t).$$

Using this fact, and the fact that each $M_\alpha(\cdot)$ is a martingale, the same argument as in the proof of Lemma 3.1 in [19] shows that

$$\langle M \rangle(t) = \sum_{\alpha \in \mathcal{T}} \Delta_\alpha^2 \bar{H}_\alpha(t).$$

We have the following deterministic upper bound on the predictable quadratic variation (quadratic characteristic):

$$\langle M \rangle(t) \leq [C_3 r] \left(\frac{C_4}{r} \right)^2 t, \quad (48)$$

where $C_3 r$ is a uniform upper bound on the total rate of potential transitions from any state, C_4/r is a uniform upper bound on Δ_α , and C_3 and C_4 are positive constants that do not depend on r .

By Doob's inequality (see e.g., Theorem 1.9.1 in [14]), for any $\delta > 0$,

$$\mathbb{P} \left\{ \max_{0 \leq \xi \leq t} |M(\xi)| \geq \delta \right\} \leq \frac{\mathbb{E}[\langle M \rangle(t)]}{\delta^2}. \quad (49)$$

In particular, for $\delta = \eta r^{3s-3-\varepsilon}$ with some $\eta \in (0, 1/4)$, and $t = T = r^{s-1}$, we have

$$\mathbb{P} \left\{ \max_{0 \leq \xi \leq T} |M(\xi)| \geq \eta r^{3s-3-\varepsilon} \right\} \leq \frac{[C_3 \lambda r] (C_4/r)^2 r^{s-1}}{[\eta r^{3s-3-\varepsilon}]^2} \leq C_{11} \frac{r^{4-5s+2\varepsilon}}{\eta^2}. \quad (50)$$

Our next goal is to estimate $|AL(\mathbf{x}) - \Xi(\mathbf{x})|$, where $\Xi(\mathbf{x})$ is defined in (25) and (29). Denote by $\bar{A}L(\mathbf{x})$ the ‘‘linear component’’ of $AL(\mathbf{x})$, i.e. the $AL(\mathbf{x})$ computed with function L replaced by its linearization at \mathbf{x} . More specifically,

$$\bar{A}L(\mathbf{x}) = \sum_{\alpha \in \mathcal{T}(\mathbf{x})} \beta_\alpha(\mathbf{x}_\alpha - \mathbf{x}) \cdot \nabla L(\mathbf{x}),$$

where $\nabla L(\mathbf{x})$ is the gradient of L at \mathbf{x} . Since the difference between $L(\mathbf{x}_\alpha) - L(\mathbf{x})$ and $(\mathbf{x}_\alpha - \mathbf{x}) \cdot \nabla L(\mathbf{x})$ is second order, using (42), we have

$$|AL(\mathbf{x}) - \bar{A}L(\mathbf{x})| \leq C_6 r \cdot \frac{r^{1-s}}{\log r} \cdot \left(\frac{1}{r} \right)^2 < r^{-s}, \quad (51)$$

where the second inequality is true for $r > \exp(C_6)$.

Let us now consider the relation between $\bar{A}L(\mathbf{x})$ and $\Xi(\mathbf{x})$. Expression for $\bar{A}L(\mathbf{x})$ can be rewritten as

$$\bar{A}L(\mathbf{x}) = \sum_{(\mathbf{k}, i) \in \mathcal{M}} \left[\frac{\partial}{\partial x_{\mathbf{k}}} L^{(a)}(\mathbf{x}) - \frac{\partial}{\partial x_{\mathbf{k}-e_i}} L^{(a)}(\mathbf{x}) \right] [v_{\mathbf{k}, i} - \tilde{w}_{\mathbf{k}, i}], \quad (52)$$

where $v_{\mathbf{k}, i}$ are obtained from making the following modifications to $\tilde{v}_{\mathbf{k}, i}$ in (29) and (27). For each i , replace $\tilde{\lambda}_i$ by λ_i (since λ_i are the *actual* fluid-scaled arrival rates), which changes the expression of $\tilde{v}_{\mathbf{k}, i}$ in (27) accordingly. Next, $x_{(i)}$ is defined as in (22), but $x_0 = a$ is replaced by the actual fluid-scaled number of zero-servers under GRAND(Z^p), i.e., $x_0 = Z^p/r$, which further changes the value of $\tilde{v}_{\mathbf{k}, i}$ in (27). We keep

the convention $(\partial/\partial x_{\mathbf{0}})L(\mathbf{x}) = 0$, for any \mathbf{x} . This completes the description of the modifications. We see that $\bar{A}L(\mathbf{x})$ is exactly equal to $\Xi(\mathbf{x})$ with each $\tilde{v}_{\mathbf{k},i}$ replaced by the corresponding $v_{\mathbf{k},i}$.

We can now estimate $|\bar{A}L(\mathbf{x}) - \Xi(\mathbf{x})|$. We know from (39) that

$$\left| \tilde{\lambda}_i - \lambda_i \right| \leq r^{-1/2+\varepsilon}.$$

Furthermore, from (38), we have

$$\left| \frac{Z^p}{r^p} - 1 \right| \leq 2r^{-1/2+\varepsilon}.$$

Therefore, the modified (and *actual*) $x_{(i)}$ with $x_{\mathbf{0}} = Z^p/r$, and the value of $x_{(i)}$ in (22), assuming $x_{\mathbf{0}} = a = r^{p-1}$, are different by at most $C'r^{p-1}r^{-1/2+\varepsilon}$. Summarizing these observations, we conclude that vectors $\tilde{\mathbf{v}} = \{\tilde{v}_{\mathbf{k},i}\}$ and $\mathbf{v} = \{v_{\mathbf{k},i}\}$ must be close, specifically

$$\|\tilde{\mathbf{v}} - \mathbf{v}\| \leq C_7 r^{-1/2+\varepsilon}. \quad (53)$$

Using the expressions (29) for $\Xi(\mathbf{x})$ and (52) for $\bar{A}L(\mathbf{x})$, we have

$$|\Xi(\mathbf{x}) - \bar{A}L(\mathbf{x})| = \sum_{(\mathbf{k},i) \in \mathcal{M}} \left[\frac{\partial}{\partial x_{\mathbf{k}}} L^{(a)}(\mathbf{x}) - \frac{\partial}{\partial x_{\mathbf{k}-\mathbf{e}_i}} L^{(a)}(\mathbf{x}) \right] [v_{\mathbf{k},i} - \tilde{v}_{\mathbf{k},i}] \quad (54)$$

$$\leq C_8 \|\nabla L^{(a)}(\mathbf{x})\| \|\tilde{\mathbf{v}} - \mathbf{v}\| \quad (55)$$

$$\leq C_9 r^{-1/2+\varepsilon}, \quad (56)$$

where the last inequality follows from (41) and (53).

Summing (51) and (56), recalling that $s > 7/8$, and rechoosing C_9 to be larger, we obtain the estimate

$$|AL(\mathbf{x}) - \Xi(\mathbf{x})| \leq |AL(\mathbf{x}) - \bar{A}L(\mathbf{x})| + |\bar{A}L(\mathbf{x}) - \Xi(\mathbf{x})| \leq C_9 r^{-1/2+\varepsilon}.$$

Then,

$$\int_0^T |AL(\mathbf{x}(t)) - \Xi(\mathbf{x}(t))| dt \leq C_9 r^{s-3/2+\varepsilon}. \quad (57)$$

By the remark after (25), we have $\Xi(\mathbf{x}) \leq 0$. Furthermore, if $|\chi_{\mathbf{k},\mathbf{k}',i}| \geq \delta_1 > 0$ for at least one pair of edges, then for $\varepsilon > 0$ and some $C = C(\delta_1) > 0$,

$$\Xi(\mathbf{x}) \leq -C(1/\log r)[r^{s-1}]^2 < -r^{2s-2-\varepsilon}.$$

We conclude that, always,

$$\int_0^t \Xi(\mathbf{x}(\xi)) d\xi \leq 0, \quad \forall t \leq T, \quad (58)$$

and if $|\chi_{\mathbf{k},\mathbf{k}',i}| \geq \delta_1 > 0$ for at least one pair of edges, in the entire interval $[0, T]$, then

$$\int_0^T \Xi(\mathbf{x}(t)) dt \leq -C(1/\log r)[r^{s-1}]^2 r^{s-1} < -r^{3s-3-\varepsilon}. \quad (59)$$

Now we specify the choice of the constant $\varepsilon > 0$: it has to be small enough to satisfy

$$3s - 3 - \varepsilon > s - 3/2 + \varepsilon \quad (\text{or } \varepsilon < s - 3/4), \quad \text{and} \quad (60)$$

$$(-3s + 3 + \varepsilon) + (4 - 5s + 2\varepsilon) < 0 \quad (\text{or } \varepsilon < (8s - 7)/3). \quad (61)$$

Since $s > 7/8$, we have both $s - 3/4 > 0$ and $(8s - 7)/3 > 0$. Thus, conditions (60) and (61) are well-defined. Now, compare $\eta r^{3s-3-\varepsilon}$, the term on the RHS of the event in bracket in (50), and $C_9 r^{s-3/2+\varepsilon}$, the term on

the RHS of the error estimate in (57). By condition (60), we have $\eta r^{3s-3-\varepsilon} > C_9 r^{s-3/2+\varepsilon}$ for all sufficiently large r . Then, from (43), (50) and (57), we obtain

$$\mathbb{P} \left\{ \max_{0 \leq t \leq T} \left| F(t) - \int_0^t \Xi(\mathbf{x}(t)) dt \right| \geq 2\eta r^{3s-3-\varepsilon} \right\} \leq C_{11} r^{4-5s+2\varepsilon} / \eta^2. \quad (62)$$

Let us remark that to obtain (62), we only need ε to be positive and satisfy condition (60). The condition (61) is needed for the rest of the proof of Theorem 4 in the next subsection.

6.4 Completing the Proof of Theorem 4

We start with a comment about the probability space constructions. The construction that we will use in this section does *not* have to be same as those used in Subsection 6.3 or Appendices A-B. The probability spaces used in different parts of the paper can be different, as long as we use them to obtain statements in the form of probability estimates, which is the case here. Obviously, the probability estimates are valid regardless of the probability space construction used to derive them.

Let the constants $\varepsilon > 0$, $\eta > 0$ and $c > 0$ be those chosen above in Section 6.3. More specifically, $\varepsilon > 0$ satisfies conditions (60) and (61), $\eta \in (0, 1/4)$, and $c > 0$ is such that Lemma 11 holds. We consider the *artificial* process on $C' r^{-3s+3+\varepsilon}$ consecutive r^{s-1} -long intervals, starting time 0. Its initial condition follows the stationary distribution of the original process. We use the following convention. If the initial state satisfies (38)-(40), then the process is according to the definition of the artificial process; otherwise, the system state is “frozen” and equal to the initial state at all times.

Let us apply the estimate (62) to each of the r^{s-1} -long intervals. (When (62) is applied to a given interval, the time is shifted so that $t = 0$ is the beginning of that interval.) Then, by (62) and a simple union bound, we see that the probability that the event in the brackets in the LHS of (62) holds for at least one of the intervals is upper bounded by

$$C_{11} C' r^{-3s+3+\varepsilon} r^{4-5s+2\varepsilon} / \eta^2 = C_{11} C' r^{7-8s+3\varepsilon} / \eta^2.$$

Since ε satisfies (61), $C_{11} C' r^{7-8s+3\varepsilon} / \eta^2 \rightarrow 0$ as $r \rightarrow \infty$. Consider a subsequence of r , increasing fast enough, e.g. $r = r(n) = e^n$, so that the sum of these probabilities is finite. Then, by Borel-Cantelli lemma, w.p.1, for all large r , the condition

$$\max_{0 \leq t \leq T} \left| F(t) - \int_0^t \tilde{A}L(\mathbf{x}(t)) dt \right| < 2\eta r^{3s-3-\varepsilon} \quad (63)$$

holds simultaneously for all $C' r^{-3s+3+\varepsilon}$ intervals. Furthermore, we have (58)-(59) for all these subintervals simultaneously.

Now consider our original system in a stationary regime for each r . (That is, its initial state coincides with that of the artificial process defined above.) We can and do construct both the artificial and original process on a common probability space, so that if the (common) initial state satisfies (38)-(40), the process trajectories coincide until the first time when the original process violates (38)-(40). By Lemmas 10 and 11, and Borel-Cantelli lemma, we can choose a further subsequence of r , along which w.p.1, for all large r , conditions (38)-(40) hold on $C' r^{-3s+3+\varepsilon}$ consecutive r^{s-1} -long intervals, and therefore the artificial process and the actual process coincide. Therefore, w.p.1, for all large r , $F(t)$ is the actual increment of the objective for all subintervals simultaneously, and we also have (58)-(59) for all subintervals simultaneously. Given this, we will apply Lemma 9 as follows.

We fix an arbitrarily small constant $\gamma > 0$. Then (w.p.1, for all large r) the following occurs. If at the beginning of a subinterval, $\Delta L \equiv L^{(a)} - L^* \geq \gamma$, then either condition $\Delta L \leq \gamma$ is “hit” within the subinterval, or, at the end of the subinterval, ΔL is smaller by at least $\delta_2 = \delta_2(r) = r^{3s-3-\varepsilon_2} / 2 > 0$. The latter is from (59), (63), and the condition that $\eta < 1/4$, because $\Delta L \geq \gamma$ implies that $|\chi_{\mathbf{k}, \mathbf{k}', i}| \geq \delta = \delta(\gamma) > 0$; here δ

does not depend on r , just like γ . In addition, if at the beginning of a subinterval or any other point in it $\Delta L \leq \gamma$, then at the end of the subinterval $\Delta L < 2\gamma$. Summarizing these statements, we obtain that (w.p.1, for all large r) at the end of the last subinterval, $\Delta L < 2\gamma$. This completes the proof of Theorem 4. \square

7 Discussion

This paper continues the line of work, originated in [23], which shows that, surprisingly, extremely simple dynamic placement (packing) algorithms, such as GRAND, can be asymptotically optimal. We show that GRAND(Z^p) algorithm, which is as simple as GRAND(aZ) in [23], is asymptotically optimal in a stronger sense than GRAND(aZ). The analysis of GRAND(Z^p) is substantially more involved technically than that of GRAND(aZ), because it cannot be reduced to the analysis of fluid limits and/or local fluid limits.

We believe that our main result, Theorem 4, holds for any $p \in (0, 1)$, without the additional condition (9) that requires p to be *sufficiently close* to 1. This additional condition is needed for our technical approach to work, and is probably just technical. Removing or relaxing the additional condition on p is interesting and important from both practical and methodological point of view, and may be a subject of future work.

References

- [1] N. Bansal, A. Caprara, M. Sviridenko. A New Approximation Method for Set Covering Problems, with Applications to Multidimensional Bin Packing. *SIAM J. Comput.*, 2009, Vol.39, No.4, pp.1256-1278.
- [2] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- [3] C. Chekuri and S. Khanna. On multidimensional packing problems. *SIAM Journal on Computing*, vol. 33, no. 4, pp. 837–851, 2004.
- [4] E. G. Coffman Jr., J. Csirik, G. Galambos, S. Martello and D. Vigo. Bin packing approximation algorithms: Survey and classification. *Handbook of Combinatorial Optimization*, ed. by P. M. Pardalos et al. (Springer, New York, 2013), pp. 455–531.
- [5] E. G. Coffman Jr., M. R. Garey and D. S. Johnson. Dynamic bin packing. *SIAM Journal on Computing*, 1983, Vol. 12, No. 2, pp. 227–258.
- [6] E. G. Coffman Jr. and A.L. Stolyar. Bandwidth Packing. *Algorithmica*, 2001, Vol. 29, pp. 70–88.
- [7] J. Csirik, D. S. Johnson, C. Kenyon, J. B. Orlin, P. W. Shor, and R. R. Weber. On the Sum-of-Squares Algorithm for Bin Packing. *J.ACM*, 2006, Vol.53, pp.1-65.
- [8] D. Gamarnik. Stochastic Bandwidth Packing Process: Stability Conditions via Lyapunov Function Technique. *Queueing Systems*, 2004, Vol.48, pp.339-363.
- [9] J. Ghaderi, Y. Zhong and R. Srikant. Asymptotic optimality of Best-Fit for stochastic bin packing. *MAMA workshop in conjunction with ACM Sigmetrics*, 2014.
- [10] A. Gulati, A. Holler, M. Ji, G. Shanmuganathan, C. Waldspurger, X. Zhu. VMware Distributed Resource Management: Design, Implementation and Lessons Learned. *VMware Technical Journal*, 2012, Vol.1, No.1, pp. 45-64. <http://labs.vmware.com/publications/vmware-technical-journal>
- [11] Y. Guo, A. L. Stolyar, A. Walid. Shadow-routing based dynamic algorithms for Virtual Machine placement in a network cloud. *INFOCOM-2013*. <http://ect.bell-labs.com/who/stolyar/publications/gpd-vm-paper-inf.pdf>
- [12] V. Gupta, A. Radovanovic. Online Stochastic Bin Packing. Preprint, 2012. <http://arxiv.org/abs/1211.2687>

- [13] J. W. Jiang, T. Lan, S. Ha, M. Chen, M. Chiang. Joint VM Placement and Routing for Data Center Traffic Engineering. *INFOCOM-2012*.
- [14] R. Sh. Liptser and A. N. Shiryaev. *Theory of Martingales*. Kluwer Academic Publishers, 1989.
- [15] S. T. Maguluri, R. Srikant, L.Ying. Stochastic Models of Load Balancing and Scheduling in Cloud Computing Clusters. *INFOCOM-2012*.
- [16] S. T. Maguluri, R. Srikant. Scheduling Jobs with Unknown Duration in Clouds. *INFOCOM-2013*.
- [17] A. Mandelbaum, W. A. Massey and M. I. Reiman. Strong Approximations for Markovian Service Networks. *Queueing Systems*, Vol. 30, 1998, pp. 149-201.
- [18] M. Mitzenmacher and E. Upfal. *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, 2005.
- [19] G. Pang, R. Talreja, W. Whitt. Martingale proofs of many-server heavy-traffic limits for Markovian queues. *Probability Surveys*, Vol. 4, 2007, pp. 193-267.
- [20] A. L. Stolyar, T. Tezcan. Shadow routing based control of flexible multi-server pools in overload. *Operations Research*, 2011, Vol.59, No.6, pp.1427-1444.
- [21] A. L. Stolyar. An infinite server system with general packing constraints. *Operations Research*, 2013, Vol.61, No.5, pp. 1200-1217.
- [22] A. L. Stolyar, Y. Zhong. A large-scale service system with packing constraints: Minimizing the number of occupied servers. *SIGMETRICS-2013*. <http://arxiv.org/abs/1212.0875>
- [23] A. L. Stolyar, Y. Zhong. Asymptotic optimality of a greedy randomized algorithm in a large-scale service system with general packing constraints. *Queueing Systems*, 2015, Vol.79, No.2, pp. 117-143. DOI 10.1007/s11134-014-9414-x
- [24] A. L. Stolyar. Large-scale heterogeneous service systems with general packing constraints. Preprint, Aug. 2015. Submitted. <http://arxiv.org/abs/1508.07512>

Appendix

A Proof of Lemma 10

A.1 Preliminaries

In this subsection, we recall some basic concentration inequalities for Poisson random variables. We then describe a convention that we will follow in this section and Appendix B.

Proposition 12 (Theorem 5.4, [18]). *Let V be a Poisson random variable with mean ν . If $x > \nu$, then*

$$\mathbb{P}(V \geq x) \leq \frac{e^{-\nu}(e\nu)^x}{x^x}; \quad (64)$$

and if $x < \nu$, then

$$\mathbb{P}(V \leq x) \leq \frac{e^{-\nu}(e\nu)^x}{x^x}. \quad (65)$$

A consequence of Proposition 12 is the following concentration bounds, which will be used extensively for the proofs of Lemmas 10 and 11.

Corollary 13. *Let V be a Poisson random variable with mean ν , and let $w \in [0, \nu]$. Then,*

$$\mathbb{P}(V \geq \nu + w) \leq e^{-\frac{w^2}{4\nu}}, \text{ and} \quad (66)$$

$$\mathbb{P}(V \leq \nu - w) \leq e^{-\frac{w^2}{4\nu}}. \quad (67)$$

Proof sketch. Let $w \in [0, \nu]$. Then, by (64),

$$\mathbb{P}(V \geq \nu + w) \leq \frac{e^{-\nu}(e\nu)^{\nu+w}}{(\nu+w)^{\nu+w}} = \exp \left[w + (\nu+w) \log \frac{\nu}{\nu+w} \right].$$

To establish (66), it suffices to prove that for all $w \in [0, \nu]$,

$$w + (\nu+w) \log \frac{\nu}{\nu+w} \leq -\frac{w^2}{4\nu}.$$

The inequality can be established by observing that when $w = 0$, LHS = RHS = 0, and that for $w \in [0, \nu]$,

$$\frac{d}{dw} \left[w + (\nu+w) \log \frac{\nu}{\nu+w} \right] = \log \frac{\nu}{\nu+w} \leq -\frac{w}{\nu+w} \leq -\frac{w}{2\nu} = \frac{d}{dw} \left[-\frac{w^2}{4\nu} \right].$$

(67) can be established in a similar way. We omit further details. □

To prove Lemmas 10 and 11, we will often consider probability bounds for a sequence of events, indexed by r . To simplify exposition, we use the following notation. For a sequence of constants $\{C^{(r)}\}$ with $C^{(r)} \in [0, 1]$ for each r , we write

$$C^{(r)} \leq e^{-\text{poly}(r)} \left(\text{respectively, } C^{(r)} \geq 1 - e^{-\text{poly}(r)} \right), \quad (68)$$

if there exists a positive constant γ such that for all sufficiently large r ,

$$C^{(r)} \leq e^{-r^\gamma} \left(\text{respectively, } C^{(r)} \geq 1 - e^{-r^\gamma} \right).$$

It is useful to think of $C^{(r)}$ as probabilities of events indexed by r . Note that if $\{C_1^{(r)}\}_r$ and $\{C_2^{(r)}\}_r$ are two sequences with $C_1^{(r)} \leq e^{-\text{poly}(r)}$ and $C_2^{(r)} \leq e^{-\text{poly}(r)}$, then we also have

$$C_1^{(r)} + C_2^{(r)} \leq e^{-\text{poly}(r)}.$$

Similarly, if $C_1^{(r)} \geq 1 - e^{-\text{poly}(r)}$ and $C_2^{(r)} \geq 1 - e^{-\text{poly}(r)}$, then

$$C_1^{(r)} + C_2^{(r)} \geq 1 - e^{-\text{poly}(r)}.$$

Finally, we also often use the expression ‘‘with probability $1 - e^{-\text{poly}(r)}$ ’’ to mean that the probability is at least $1 - e^{-r^\gamma}$, for some $\gamma > 0$ and for all sufficiently large r .

A.2 Proof of Lemma 10

Our general strategy for establishing Lemma 10 is to divide the interval $[0, r^\alpha]$ into shorter sub-intervals of length $r^{-1/2}$, and show that (39) holds with high probability over each sub-interval. More specifically, let us define events

$$E_{j,i}^r = \left\{ \sup_{t \in [(j-1)r^{-\frac{1}{2}}, jr^{-\frac{1}{2}}]} |Y_i^r(t) - \rho_i r| \geq \frac{r^{1/2+\varepsilon}}{I} \right\}, \quad (69)$$

for $i \in \mathcal{I}$, $j \in \{1, 2, \dots, r^{\alpha+\frac{1}{2}}\}^1$, and for each r . Note that because for each r , the system is in the stationary regime, we have that for any $i \in \mathcal{I}$,

$$\mathbb{P}(E_{1,i}^r) = \mathbb{P}(E_{2,i}^r) = \dots = \mathbb{P}(E_{r^{\alpha+1/2},i}^r). \quad (70)$$

Thus, if we can show that

$$\mathbb{P}(E_{1,i}^r) \leq e^{-\text{poly}(r)}, \forall i \in \mathcal{I}, \quad (71)$$

then by (70) and (71),

$$\mathbb{P}(\cup_{1 \leq j \leq r^{\alpha+1/2}, i \in \mathcal{I}} E_{j,i}^r) \leq e^{-\text{poly}(r)} \rightarrow 0 \quad \text{as } r \rightarrow \infty,$$

which establishes Lemma 10 immediately.

It is now left to show that $\mathbb{P}(E_{1,i}^r) \leq e^{-\text{poly}(r)}$ for each $i \in \mathcal{I}$.

Lemma 14. *Let $i \in \mathcal{I}$, and consider the system indexed by r in the stationary regime. Then, for sufficiently large r ,*

$$\mathbb{P}(E_{1,i}^r) = \mathbb{P}\left(\sup_{t \in [0, r^{-1/2}]} |Y_i^r(t) - \rho_i r| \geq \frac{r^{1/2+\varepsilon}}{I}\right) \leq e^{-\text{poly}(r)}. \quad (72)$$

Proof. For notational convenience, we drop the subscript i . The proof of Lemma 14 consists of establishing (a) a probability tail bound of $|Y^r(0) - \rho r|$; and (b) a probability tail bound of $\sup_{t \in [0, r^{-1/2}]} (Y^r(t) - \rho r)$; and (c) a probability tail bound of $\inf_{t \in [0, r^{-1/2}]} (Y^r(t) - \rho r)$.

(a) We first show that

$$\mathbb{P}\left(|Y^r(0) - \rho r| \geq r^{1/2+\varepsilon/3}\right) \leq e^{-\text{poly}(r)}. \quad (73)$$

Since $Y^r(0)$ is a Poisson random variable with mean ρr , we apply the concentration bounds (66) and (67). First by (66), we have that for sufficiently large r , $r^{1/2+\varepsilon/3} \leq \rho r$, and

$$\mathbb{P}\left(Y^r(0) - \rho r \geq r^{1/2+\varepsilon/3}\right) \leq \exp\left(-\frac{r^{1+2\varepsilon/3}}{4\rho r}\right) = \exp\left(-\frac{r^{2\varepsilon/3}}{4\rho}\right).$$

¹We treat $r^{\alpha+\frac{1}{2}}$ as if it were guaranteed to be an integer. Similarly, in the sequel, whenever Cr^β is used as an index for some constants $\beta > 0$ and $C > 0$, we treat it as an integer. Rounding them up or down to a nearest integer would overburden our notation, but would not affect our order-of-magnitude estimates.

Similarly, for sufficiently large r ,

$$\mathbb{P}\left(Y^r(0) - \rho r \leq -r^{1/2+\varepsilon/3}\right) \leq \exp\left(-\frac{r^{2\varepsilon/3}}{4\rho}\right).$$

Thus, by a simple union bound, for sufficiently large r ,

$$\mathbb{P}\left(|Y^r(0) - \rho r| \geq r^{1/2+\varepsilon/3}\right) \leq 2 \exp\left(-\frac{r^{2\varepsilon/3}}{4\rho}\right) \leq e^{-\text{poly}(r)}. \quad (74)$$

This completes part (a).

(b) We then show that

$$\mathbb{P}\left(\sup_{t \in [0, r^{-1/2}]} (Y^r(t) - \rho r) \geq r^{1/2+2\varepsilon/3}\right) \leq e^{-\text{poly}(r)}. \quad (75)$$

To establish (75), we use the following representation of the process $Y^r(\cdot)$. Note that the process $Y^r(\cdot)$ describes the steady-state evolution of a $M/M/\infty$ queueing system, so we can represent $Y^r(\cdot)$ as (see e.g., [19])

$$Y^r(t) = Y^r(0) + \Pi(\lambda r t) - \tilde{\Pi}\left(\int_0^t \mu Y^r(s) ds\right), \quad (76)$$

where $Y^r(0)$ is a Poisson random variable with mean ρr , and $\Pi(\cdot)$ and $\tilde{\Pi}(\cdot)$ are independent unit-rate Poisson processes that are also independent from $Y^r(0)$.

By (76), we have that w.p.1, for all $t \in [0, r^{-1/2}]$,

$$Y^r(t) \leq Y^r(0) + \Pi(\lambda r t) \leq Y^r(0) + \Pi\left(\lambda r \cdot r^{-1/2}\right) = Y^r(0) + \Pi\left(\lambda r^{1/2}\right).$$

Using (66) and the fact that $\Pi\left(\lambda r^{1/2}\right)$ is Poisson with mean $\lambda r^{1/2}$, we have that

$$\mathbb{P}\left(\Pi\left(\lambda r^{1/2}\right) \geq 2\lambda r^{1/2}\right) \leq \exp\left(-\frac{\lambda r^{1/2}}{4}\right) \leq e^{-\text{poly}(r)}. \quad (77)$$

For sufficiently large r , if $Y^r(0) + \Pi\left(\lambda r^{1/2}\right) \geq \rho_i r + r^{1/2+2\varepsilon/3}$, then we must have $Y^r(0) \geq \rho_i r + r^{1/2+\varepsilon/3}$ or $\Pi\left(\lambda r^{1/2}\right) \geq 2\lambda r^{1/2}$. Thus,

$$\begin{aligned} \mathbb{P}\left(Y^r(0) + \Pi\left(\lambda r^{1/2}\right) \geq \rho_i r + r^{1/2+2\varepsilon/3}\right) &\leq \mathbb{P}\left(Y^r(0) \geq \rho_i r + r^{1/2+\varepsilon/3} \text{ or } \Pi\left(\lambda r^{1/2}\right) \geq 2\lambda r^{1/2}\right) \\ &\leq \mathbb{P}\left(Y^r(0) \geq \rho_i r + r^{1/2+\varepsilon/3}\right) \\ &\quad + \mathbb{P}\left(\Pi\left(\lambda r^{1/2}\right) \geq 2\lambda r^{1/2}\right) \\ &\leq e^{-\text{poly}(r)}, \end{aligned}$$

where the last inequality follows from (73) and (77).

Since with probability 1, for all $t \in [0, r^{-1/2}]$,

$$Y^r(t) \leq Y^r(0) + \Pi\left(\lambda r^{1/2}\right),$$

it immediately follows that (75) holds. This completes part (b).

(c) We now show that

$$\mathbb{P}\left(\inf_{t \in [0, r^{-1/2}]} (Y^r(t) - \rho r) \leq -r^{1/2+2\varepsilon/3}\right) \leq e^{-\text{poly}(r)}. \quad (78)$$

Using the representation (76), we have that with probability 1, for all $t \in [0, r^{-1/2}]$,

$$Y^r(t) \geq Y^r(0) - \tilde{\Pi} \left(\int_0^t \mu Y^r(s) ds \right) \geq Y^r(0) - \tilde{\Pi} \left(\int_0^{r^{-1/2}} \mu Y^r(s) ds \right).$$

Consider events F^r and G^r defined to be

$$F^r = \left\{ \sup_{t \in [0, r^{-1/2}]} Y^r(t) < 2\rho r \right\}, \quad \text{and} \quad G^r = \left\{ \tilde{\Pi} \left(r^{1/2+\varepsilon/3} \right) < 2r^{1/2+\varepsilon/3} \right\}.$$

For sufficiently large r , we have that under event F^r ,

$$\int_0^{r^{-1/2}} \mu Y^r(s) ds < \mu r^{-1/2} \cdot (2\rho r) \leq r^{1/2+\varepsilon/3}.$$

Thus, for sufficiently large r , we have that under the event $F^r \cap G^r$,

$$\tilde{\Pi} \left(\int_0^{r^{-1/2}} \mu Y^r(s) ds \right) \leq \tilde{\Pi} \left(r^{1/2+\varepsilon/3} \right) < 2r^{1/2+\varepsilon/3}.$$

By (75), we know that $\mathbb{P}(F^r) \geq 1 - e^{-\text{poly}(r)}$, and using (66), we can easily show that $\mathbb{P}(G^r) \geq 1 - e^{-\text{poly}(r)}$. Thus, we have $\mathbb{P}(F^r \cap G^r) \geq 1 - e^{-\text{poly}(r)}$, and by considering the complement, we have

$$\mathbb{P} \left(\tilde{\Pi} \left(\int_0^{r^{-1/2}} \mu Y^r(s) ds \right) \geq 2r^{1/2+\varepsilon/3} \right) \leq e^{-\text{poly}(r)}. \quad (79)$$

For sufficiently large r , if $Y^r(0) - \tilde{\Pi} \left(\int_0^{r^{-1/2}} \mu Y^r(s) ds \right) \leq \rho_i r - r^{1/2+2\varepsilon/3}$, then we have $Y^r(0) \leq \rho_i r - r^{1/2+\varepsilon/3}$ or $\tilde{\Pi} \left(\int_0^{r^{-1/2}} \mu Y^r(s) ds \right) \geq 2r^{1/2+\varepsilon/3}$. Thus, similar to the argument in part (b), we have that

$$\begin{aligned} \mathbb{P} \left(Y^r(0) - \tilde{\Pi} \left(\int_0^{r^{-1/2}} \mu Y^r(s) ds \right) \leq \rho_i r - r^{1/2+2\varepsilon/3} \right) &\leq \mathbb{P} \left(Y^r(0) \leq \rho_i r - r^{1/2+\varepsilon/3} \right) \\ &\quad + \mathbb{P} \left(\tilde{\Pi} \left(\int_0^{r^{-1/2}} \mu Y^r(s) ds \right) \geq 2r^{1/2+\varepsilon/3} \right) \\ &\leq e^{-\text{poly}(r)}, \end{aligned}$$

where the last inequality follows from (73) and (79). This establishes (78), and completes part (c).

Combining (73), (75) and (78), we have

$$\mathbb{P} \left(\sup_{t \in [0, r^{-1/2}]} |Y^r(t) - \rho r| \geq r^{1/2+2\varepsilon/3} \right) \leq e^{-\text{poly}(r)}.$$

Since for sufficiently large r , $\frac{r^{1/2+\varepsilon}}{I} \geq r^{1/2+2\varepsilon/3}$, we also have

$$\mathbb{P} \left(\sup_{t \in [0, r^{-1/2}]} |Y^r(t) - \rho r| \geq \frac{r^{1/2+\varepsilon}}{I} \right) \leq \mathbb{P} \left(\sup_{t \in [0, r^{-1/2}]} |Y^r(t) - \rho r| \geq r^{1/2+2\varepsilon/3} \right) \leq e^{-\text{poly}(r)}.$$

This establishes Lemma 14. □

B Proof of Lemma 11

For each $\mathbf{k} \in \bar{\mathcal{K}}$, define constant $s(\mathbf{k}) = 1 - (1 + \sum_{i \in \mathcal{I}} k_i)(1 - p)$. To establish Lemma 11, we consider instead the condition

$$X_{\mathbf{k}}^r(t) \geq cr^{s(\mathbf{k})} \quad (80)$$

for each $\mathbf{k} \in \bar{\mathcal{K}}$, at any time $t \geq 0$, and prove the following stronger result.

Lemma 15. *Let $\alpha > 0$. Consider our sequence of systems in the stationary regime, under the GRAND(Z^p) algorithm. Then for each $\mathbf{k} \in \bar{\mathcal{K}}$, there exists a constant $c > 0$, such that as $r \rightarrow \infty$,*

$$\mathbb{P}(\text{Condition (80) holds for all } t \in [0, r^\alpha]) \geq 1 - e^{-\text{poly}(r)}. \quad (81)$$

To prove Lemma 15, we use the following construction of the underlying probability space. For each $(\mathbf{k}, i) \in \mathcal{M}$, let $\Pi_{(\mathbf{k}, i)}(\cdot)$ and $\tilde{\Pi}_{(\mathbf{k}, i)}(\cdot)$ be unit-rate Poisson processes. Furthermore, all $\Pi_{(\mathbf{k}, i)}(\cdot)$ and $\tilde{\Pi}_{(\mathbf{k}, i)}(\cdot)$ are independent. We will use $\Pi_{(\mathbf{k}, i)}(\cdot)$ as the driving processes for customer arrivals, and $\tilde{\Pi}_{(\mathbf{k}, i)}(\cdot)$ to drive departures. Furthermore, processes $\Pi_{(\mathbf{k}, i)}(\cdot)$ and $\tilde{\Pi}_{(\mathbf{k}, i)}(\cdot)$ are all independent from the initial random system state $\mathbf{X}^r(0)$.

More specifically, let $D_{(\mathbf{k}, i)}^r(t)$ denote the total number of departures along the edge (\mathbf{k}, i) in $[0, t]$. Then,

$$D_{(\mathbf{k}, i)}^r(t) = \tilde{\Pi}_{(\mathbf{k}, i)} \left(\int_0^t X_{\mathbf{k}}^r(u) k_i \mu_i du \right). \quad (82)$$

Similarly, let $A_{(\mathbf{k}, i)}^r(t)$ denote the total number of arrivals along the edge (\mathbf{k}, i) in $[0, t]$. Under the GRAND algorithm, a type- i customer that arrives at time s is placed in a server of configuration $\mathbf{k} - \mathbf{e}_i$ with probability $X_{\mathbf{k} - \mathbf{e}_i}^r(s) / X_{(\mathbf{k}, i)}^r(s)$. Then, we can write

$$A_{(\mathbf{k}, i)}^r(t) = \Pi_{(\mathbf{k}, i)} \left(\int_0^t \lambda_i r \cdot \frac{X_{\mathbf{k} - \mathbf{e}_i}^r(u)}{X_{(\mathbf{k}, i)}^r(u)} du \right). \quad (83)$$

Thus, for each $\mathbf{k} \in \mathcal{K}$, we can write

$$\begin{aligned} X_{\mathbf{k}}^r(t) - X_{\mathbf{k}}^r(0) &= \left[\sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} A_{(\mathbf{k}, i)}^r(t) + \sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} D_{(\mathbf{k} + \mathbf{e}_i, i)}^r(t) \right] \\ &\quad - \left[\sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} A_{(\mathbf{k} + \mathbf{e}_i, i)}^r(t) + \sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} D_{(\mathbf{k}, i)}^r(t) \right] \end{aligned} \quad (84)$$

$$\begin{aligned} &\geq \left[\sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} A_{(\mathbf{k}, i)}^r(t) \right] - \left[\sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} A_{(\mathbf{k} + \mathbf{e}_i, i)}^r(t) + \sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} D_{(\mathbf{k}, i)}^r(t) \right] \\ &= S_1(t) - S_2(t) - S_3(t), \end{aligned} \quad (85)$$

where

$$S_1(t) = \sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} A_{(\mathbf{k}, i)}^r(t); \quad (86)$$

$$S_2(t) = \sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} A_{(\mathbf{k} + \mathbf{e}_i, i)}^r(t); \quad (87)$$

$$S_3(t) = \sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} D_{(\mathbf{k}, i)}^r(t). \quad (88)$$

The equality (84) follows by accounting for all arrivals and departures that contribute to $X_{\mathbf{k}}^r(t) - X_{\mathbf{k}}^r(0)$, the change in $X_{\mathbf{k}}^r$. For example, the term $\sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} A_{(\mathbf{k}, i)}^r(t)$ accounts for arrivals along the edges $(\mathbf{k}, i) \in \mathcal{M}$, which increases the number of servers of configuration \mathbf{k} .

Proof of Lemma 15. We prove Lemma 15 by induction on $\|\mathbf{k}\|_1 = \sum_i k_i$.

Base case: $\|\mathbf{k}\|_1 = 0$. If $\|\mathbf{k}\|_1 = 0$, then $\mathbf{k} = \mathbf{0}$, and $s(\mathbf{0}) = 1 - (1 - p) = p$. Then, we show that

$$\mathbb{P}\left(X_{\mathbf{0}}^r(t) \geq \frac{1}{2}r^p, \forall t \in [0, r^\alpha]\right) \geq 1 - e^{-\text{poly}(r)}.$$

To this end, note that $X_{\mathbf{0}}^r(t) = \lceil (Z^r(t))^p \rceil$ for all t . By Lemma 10, we know that for any $\varepsilon > 0$, as $r \rightarrow \infty$,

$$\mathbb{P}\left(|Y_i^r(t) - \rho_i r| \leq \frac{r^{1/2+\varepsilon}}{I}, \forall t \in [0, r^\alpha], \forall i \in \mathcal{I}\right) \geq 1 - e^{-\text{poly}(r)}.$$

In particular, by setting $\varepsilon = 1/4$, we have

$$\mathbb{P}\left(|Y_i^r(t) - \rho_i r| \leq \frac{r^{3/4}}{I}, \forall t \in [0, r^\alpha], \forall i \in \mathcal{I}\right) \geq 1 - e^{-\text{poly}(r)}.$$

For sufficiently large r , $\frac{r}{2^{1/p}} \leq r - r^{3/4}$. Since $Z^r(t) = \sum_{i \in \mathcal{I}} Y_i^r(t)$ for all t , we have

$$\begin{aligned} \mathbb{P}\left(Z^r(t) \geq \frac{r}{2^{1/p}}, \forall t \in [0, r^\alpha]\right) &\geq \mathbb{P}\left(|Z^r(t) - r| \leq r^{3/4}, \forall t \in [0, r^\alpha]\right) \\ &\geq \mathbb{P}\left(|Y_i^r(t) - \rho_i r| \leq \frac{r^{3/4}}{I}, \forall t \in [0, r^\alpha], \forall i \in \mathcal{I}\right) \\ &\geq 1 - e^{-\text{poly}(r)}. \end{aligned}$$

This implies that

$$\mathbb{P}\left(X_{\mathbf{0}}^r(t) \geq \frac{1}{2}r^p, \forall t \in [0, r^\alpha]\right) \geq \mathbb{P}\left(Z^r(t) \geq \frac{r}{2^{1/p}}, \forall t \in [0, r^\alpha]\right) \geq 1 - e^{-\text{poly}(r)},$$

establishing the base case.

Induction step. Let $\ell > 0$. Suppose that for all \mathbf{k}' with $\|\mathbf{k}'\|_1 \leq \ell - 1$, there exists $c' > 0$ such that

$$\mathbb{P}\left(X_{\mathbf{k}'}^r(t) \geq c' r^{s(\mathbf{k}')}, \forall t \in [0, r^\alpha]\right) \geq 1 - e^{-\text{poly}(r)}.$$

Let $\mathbf{k} \in \mathcal{K}$ be a configuration with $\|\mathbf{k}\|_1 = \ell > 0$. Then, there exists $i \in \mathcal{I}$ such that $k_i \geq 1$. Let $\tilde{\mathbf{k}} = \mathbf{k} - \mathbf{e}_i$ so that $\|\tilde{\mathbf{k}}\|_1 = \ell - 1$. For notational convenience, write $s_1 = s(\tilde{\mathbf{k}})$ and $s_2 = s_1 - (1 - p)$ so that $s_2 = s(\mathbf{k})$. By the induction hypothesis, we can assume that

$$\mathbb{P}\left(X_{\tilde{\mathbf{k}}}^r(t) \geq \tilde{c} r^{s_1}, \forall t \in [0, r^\alpha]\right) \geq 1 - e^{-\text{poly}(r)} \quad (89)$$

for some $\tilde{c} > 0$. Furthermore, let E^r be the event

$$E^r = \{X_{\tilde{\mathbf{k}}}^r(t) \geq \tilde{c} r^{s_1}, \forall t \in [0, r^\alpha]\} \quad (90)$$

so that $\mathbb{P}(E^r) \geq 1 - e^{-\text{poly}(r)}$.

Consider the interval $[0, r^\alpha + 1]$, and divide it into $\frac{r^\alpha + 1}{T}$ sub-intervals of length $T = r^{s_2 - 1 - \varepsilon'}$, where $0 < \varepsilon' < s_2 - \frac{p}{2}$; namely, sub-intervals $[0, T], [T, 2T], [2T, 3T], \dots$. We claim the following.

Claim 1. There exist positive constants c and \bar{c} such that with probability $1 - e^{-\text{poly}(r)}$, for each $j \in \{1, 2, \dots, \frac{r^\alpha + 1}{T}\}$:

(i) if $X_{\tilde{\mathbf{k}}}^r((j-1)T) \leq 4cr^{s_2}$, then

$$X_{\tilde{\mathbf{k}}}^r(jT) - X_{\tilde{\mathbf{k}}}^r((j-1)T) \geq \bar{c} r^{2s_2 - p - \varepsilon'};$$

(ii) if $X_{\mathbf{k}}^r((j-1)T) \geq 2cr^{s_2}$, then

$$\inf_{t \in [(j-1)T, jT]} X_{\mathbf{k}}^r(t) \geq \frac{1}{2} X_{\mathbf{k}}^r(0).$$

Assuming the validity of Claim 1, we now complete the rest of the induction step, and defer the proof of the claim.

Suppose that Claim 1 is true. For each r , consider a sample path for which both statements (i) and (ii) of Claim 1 hold for all $j \in \{1, 2, \dots, \frac{r^\alpha+1}{T}\}$, and let j_0 be minimal such that

$$X_{\mathbf{k}}^r(j_0T) \geq 4cr^{s_2}. \quad (91)$$

Then, for sufficiently large r , $j_0 \leq \frac{1}{T}$. If $X_{\mathbf{k}}^r(0) \geq 4cr^{s_2}$, then $j_0 = 0 \leq \frac{1}{T}$. If $X_{\mathbf{k}}^r(0) < 4cr^{s_2}$, then we also have $j_0 \leq \frac{1}{T}$, since by comparing the threshold $4cr^{s_2}$ and increment $\bar{c}r^{2s_2-p-\varepsilon'}$ from $X_{\mathbf{k}}^r(jT)$ to $X_{\mathbf{k}}^r((j-1)T)$ in statement (i), we have

$$\frac{4cr^{s_2}}{\bar{c}r^{2s_2-p-\varepsilon'}} = \frac{4c}{\bar{c}} \cdot \frac{1}{r^{s_2-p-\varepsilon'}} = \frac{4c}{\bar{c}} \cdot \frac{r^{p-1}}{T} \leq \frac{1}{T}.$$

Next, for any $j \in \{j_0, j_0+1, \dots, \frac{r^\alpha+1}{T}\}$, we have

$$X_{\mathbf{k}}^r(jT) \geq 2cr^{s_2}, \quad (92)$$

which we establish by induction (note that this induction is distinct from the induction on $\|\mathbf{k}\|$ that we use to prove Lemma 15). (92) is true for $j = j_0$, by (91). For $j > j_0$, suppose by induction hypothesis that $X_{\mathbf{k}}^r((j-1)T) \geq 2cr^{s_2}$. Then, we consider two cases: if $X_{\mathbf{k}}^r((j-1)T) \geq 4cr^{s_2}$, then by statement (ii), (92) holds for j ; if $X_{\mathbf{k}}^r((j-1)T) \in [2cr^{s_2}, 4cr^{s_2})$, then by statement (i), (92) holds for j as well. This completes the induction step and the proof of (92) for all $j \in \{j_0, j_0+1, \dots, \frac{r^\alpha+1}{T}\}$. By statement (ii), it then follows that for all $t \in [1, r^\alpha+1]$,

$$X_{\mathbf{k}}^r(t) \geq cr^{s_2}. \quad (93)$$

In summary, we have established that with probability $1 - e^{-\text{poly}(r)}$, for all $t \in [1, r^\alpha+1]$,

$$X_{\mathbf{k}}^r(t) \geq cr^{s_2} = cr^{s(\mathbf{k})}.$$

By the stationarity of the processes $\mathbf{X}^r(\cdot)$, we can conclude that

$$\mathbb{P}(X_{\mathbf{k}}^r(t) \geq cr^{s_2}, \forall t \in [0, r^\alpha]) \geq 1 - e^{-\text{poly}(r)}.$$

This completes the induction step, assuming that Claim 1 holds. It is now left to prove Claim 1.

Proof of Claim 1. Let us first focus on the sub-interval $[0, T]$. Again, since the systems indexed by r are in the stationary regime, probability bounds that we derive over the sub-interval $[0, T]$ extend naturally to other sub-intervals $[jT, (j+1)T]$, $j = 1, 2, \dots$. Claim 1 is an easy consequence of the following claim.

Claim 2. There exist positive constants c and \bar{c} such that with probability $1 - e^{-\text{poly}(r)}$, the following holds.

(1) If $X_{\mathbf{k}}^r(0) \leq 4cr^{s_2}$, then

$$X_{\mathbf{k}}^r(T) - X_{\mathbf{k}}^r(0) \geq \bar{c}r^{2s_2-p-\varepsilon'}. \quad (94)$$

(2) If $X_{\mathbf{k}}^r(0) \geq 2cr^{s_2}$, then

$$\inf_{t \in [0, T]} X_{\mathbf{k}}^r(t) \geq \frac{1}{2} X_{\mathbf{k}}^r(0). \quad (95)$$

To see how Claim 2 implies Claim 1, recall that our systems are in the stationary regime. Thus, for any other sub-interval $[jT, (j+1)T]$, with the same probability $1 - e^{-\text{poly}(r)}$, statements (1) and (2) of Claim 2 hold, when we replace 0 by jT and T by $(j+1)T$. Since there are a polynomial (in r) number of such sub-intervals, Claim 1 then follows immediately. We now prove Claim 2.

Proof of Claim 2. To bound $X_{\mathbf{k}}^r(\cdot)$, we consider terms $S_1(T), S_2(T)$ and $S_3(T)$ defined in (86) – (88) separately.

(a) First, consider the term $S_1(T)$ defined in (86). Focus on $A_{(\mathbf{k}, \iota)}^r(T)$, the cumulative arrivals along the edge (\mathbf{k}, ι) , where we recall that $k_\iota \geq 1$, and $\tilde{\mathbf{k}} = \mathbf{k} - \mathbf{e}_\iota$. We have

$$A_{(\mathbf{k}, \iota)}^r(T) = \Pi_{(\mathbf{k}, \iota)} \left(\int_0^T \lambda_\iota r \frac{X_{\tilde{\mathbf{k}}}^r(u)}{X_{(\iota)}^r(u)} du \right) \leq S_1(T).$$

Let F_1^r be the event defined by

$$F_1^r = \left\{ \frac{1}{2}r \leq Z^r(u) \leq \frac{3}{2}r, \text{ for all } u \in [0, T] \right\}.$$

Then, by inspecting the proof of the base case, it is easy to see that $\mathbb{P}(F_1^r) \geq 1 - e^{-\text{poly}(r)}$. Furthermore, for sufficiently large r , under the event F_1^r , we have that for all $u \in [0, T]$,

$$X_{(\iota)}^r(u) \leq Z^r(u) + [Z^r(u)]^p \leq 2r.$$

Thus, for sufficiently large r , under the event $E^r \cap F_1^r$, where we recall the definition of E^r in (90), we have that for all $u \in [0, T]$,

$$\frac{X_{\tilde{\mathbf{k}}}^r(u)}{X_{(\iota)}^r(u)} \geq \frac{\tilde{c}r^{s_1}}{2r} = \frac{\tilde{c}}{2}r^{s_1-1},$$

from which it follows that

$$\int_0^T \lambda_\iota r \frac{X_{\tilde{\mathbf{k}}}^r(u)}{X_{(\iota)}^r(u)} du \geq \lambda_\iota r \cdot \frac{\tilde{c}}{2}r^{s_1-1} \cdot T = c_1 r^{s_1+s_2-1-\varepsilon'}.$$

where $c_1 = \frac{1}{2}\lambda_\iota\tilde{c}$. We also define event G_1^r to be

$$G_1^r = \left\{ \Pi_{(\mathbf{k}, \iota)} \left(c_1 r^{s_1+s_2-1-\varepsilon'} \right) \geq \frac{1}{2}c_1 r^{s_1+s_2-1-\varepsilon'} \right\}.$$

Then, $\mathbb{P}(G_1^r) \geq 1 - e^{-\text{poly}(r)}$, and under the event $E^r \cap F_1^r \cap G_1^r$, we have that

$$A_{(\mathbf{k}, \iota)}^r(T) \geq \frac{1}{2}c_1 r^{s_1+s_2-1-\varepsilon'}. \quad (96)$$

Thus,

$$\mathbb{P} \left(S_1(T) \geq \frac{1}{2}c_1 r^{s_1+s_2-1-\varepsilon'} \right) \geq \mathbb{P} \left(A_{(\mathbf{k}, \iota)}^r(T) \geq \frac{1}{2}c_1 r^{s_1+s_2-1-\varepsilon'} \right) \geq 1 - e^{-\text{poly}(r)}. \quad (97)$$

This completes our probability bound for $S_1(T)$ and part (a).

(b) Next, we consider the term $S_2(T)$ defined in (87). We have

$$S_2(T) = \sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} A_{(\mathbf{k} + \mathbf{e}_i, i)}^r(T) = \sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} \Pi_{(\mathbf{k} + \mathbf{e}_i, i)} \left(\int_0^T \lambda_i r \cdot \frac{X_{\mathbf{k}}^r(u)}{X_{(i)}^r(u)} du \right).$$

Instead of deriving probability bounds for $S_2(T)$, in part (b) we will only derive a probability bound for

$$\sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} \int_0^T \lambda_i r \cdot \frac{X_{\mathbf{k}}^r(u)}{X_{(i)}^r(u)} du. \quad (98)$$

The reason is twofold. First, the term (98) captures the order of magnitude of $S_2(T)$. Second, the bound that we will derive for the term (98) depends on $X_{\mathbf{k}}^r(0)$. Since we will consider two separate cases depending on the

magnitude of $X_{\mathbf{k}}^r(0)$, we will have separate probability bounds for $S_2(T)$, which we derive after considering the term $S_3(T)$. Consider $\int_0^T \lambda_i r \frac{X_{\mathbf{k}}^r(u)}{X_{(i)}^r(u)} du$, a generic summand of the term (98). By the base case and using the fact that $X_{(i)}^r(u) \geq X_{\mathbf{0}}^r(u)$ for all u , we have that

$$\mathbb{P}\left(X_{(i)}^r(u) \geq \frac{1}{2}r^p\right) \geq \mathbb{P}\left(X_{\mathbf{0}}^r(u) \geq \frac{1}{2}r^p\right) \geq 1 - e^{-\text{poly}(r)}. \quad (99)$$

Furthermore, there exists a constant $c_2 > 0$ such that

$$\mathbb{P}(X_{\mathbf{k}}^r(u) - X_{\mathbf{k}}^r(0) \leq c_2 r T, \quad \forall u \in [0, T]) \geq 1 - e^{-\text{poly}(r)}. \quad (100)$$

We establish (100) as follows. Let N^r be the total number of arrivals to and departures from the system up to time T . Then, it is clear that for all $u \in [0, T]$, $X_{\mathbf{k}}^r(u) - X_{\mathbf{k}}^r(0) \leq N^r$. We now obtain a probability bound on N^r . Since we are only interested in probability bounds, we can and will at different points of the proof use different underlying probability space constructions, as long as they produce the same – in law – system process. At this point, we will use the following, different probability space construction. We associate a driving unit-rate Poisson process $\Pi(\cdot)$ for all the arrivals to the system, and an independent unit-rate Poisson process $\tilde{\Pi}(\cdot)$ to drive all the departures. The total arrival rate at all times is r , and the total departure rate at time t is given by $\sum_i \mu_i Y_i^r(t)$. Thus, N^r has the same distribution as $\Pi(rT) + \tilde{\Pi}\left(\int_0^T \sum_i \mu_i Y_i^r(u) du\right)$. With probability $1 - e^{-\text{poly}(r)}$, $\Pi(rT) \leq 2rT$. By Lemma 10, it is easy to see that with probability $1 - e^{-\text{poly}(r)}$, for all $u \in [0, T]$, $\sum_i \mu_i Y_i^r(u) \leq \sum_i \mu_i (2\rho_i r) = 2 \sum_i \lambda_i r = c'_2 r$ for $c'_2 = 2 \sum_i \lambda_i$. This implies that with probability $1 - e^{-\text{poly}(r)}$,

$$\tilde{\Pi}\left(\int_0^T \sum_i \mu_i Y_i^r(u) du\right) \leq \tilde{\Pi}\left(\int_0^T c'_2 r du\right) = \tilde{\Pi}(c'_2 r T) \leq 2c'_2 r T.$$

Thus, with probability $1 - e^{-\text{poly}(r)}$, for all $u \in [0, T]$,

$$X_{\mathbf{k}}^r(u) - X_{\mathbf{k}}^r(0) \leq N^r \stackrel{d}{=} \Pi(rT) + \tilde{\Pi}\left(\int_0^T \sum_i \mu_i Y_i^r(u) du\right) \leq 2rT + 2c'_2 r T = c_2 r T.$$

where $c_2 = 2 + 2c'_2$. This establishes (100).

By (99) and (100), we can bound (98) as follows. With probability $1 - e^{-\text{poly}(r)}$,

$$\begin{aligned} \int_0^T \lambda_i r \frac{X_{\mathbf{k}}^r(u)}{X_{(i)}^r(u)} du &\leq \int_0^T \lambda_i r \frac{X_{\mathbf{k}}^r(0) + c_2 r T}{r^p/2} du \\ &\leq c_3 T r^{1-p} X_{\mathbf{k}}^r(0) + c_4 T^2 r^{2-p} \\ &= c_5 r^{s_2-p-\varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2-\varepsilon'}\right) \end{aligned}$$

for some positive constants c_3 , c_4 and c_5 . It then follows immediately that with probability $1 - e^{-\text{poly}(r)}$,

$$\sum_{i: \mathbf{k} + \mathbf{e}_i \in \mathcal{K}} \int_0^T \lambda_i r \frac{X_{\mathbf{k}}^r(u)}{X_{(i')}^r(u)} du \leq c_6 r^{s_2-p-\varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2-\varepsilon'}\right), \quad (101)$$

for some positive constant c_6 . This completes part (b).

(c) We now consider the term $S_3(T)$ defined in (88), which is given by

$$S_3(T) = \sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} D_{(\mathbf{k}, i)}^r(T) = \sum_{i: \mathbf{k} - \mathbf{e}_i \in \bar{\mathcal{K}}} \tilde{\Pi}_{(\mathbf{k}, i)} \left(\int_0^T X_{\mathbf{k}}^r(u) k_i \mu_i du \right).$$

Similar to part (b), we only derive a probability bound for

$$\sum_{i:\mathbf{k}-\mathbf{e}_i \in \bar{\mathcal{K}}} \int_0^T X_{\mathbf{k}}^r(u) k_i \mu_i du. \quad (102)$$

By (100), we have that with probability $1 - e^{-\text{poly}(r)}$,

$$\begin{aligned} \int_0^T k_i \mu_i X_{\mathbf{k}}^r(u) du &\leq \int_0^T k_i \mu_i (X_{\mathbf{k}}^r(0) + c_2 r T) du \\ &= c_7 T X_{\mathbf{k}}^r(0) + c_8 r T^2 \\ &= c_9 r^{s_2-1-\varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2-\varepsilon'} \right), \end{aligned}$$

and

$$\sum_{i:\mathbf{k}-\mathbf{e}_i \in \bar{\mathcal{K}}} \int_0^T k_i \mu_i X_{\mathbf{k}}^r(u) du \leq c_{10} r^{s_2-1-\varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2-\varepsilon'} \right), \quad (103)$$

for some positive constants c_7, c_8, c_9 and c_{10} . This completes part (c).

Observe that by (101) and (103), we have that with probability $1 - e^{-\text{poly}(r)}$,

$$\sum_{i:\mathbf{k}+\mathbf{e}_i \in \mathcal{K}} \int_0^T \lambda_i r \frac{X_{\mathbf{k}}^r(u)}{X_{(i)}^r(u)} du + \sum_{i:\mathbf{k}-\mathbf{e}_i \in \bar{\mathcal{K}}} \int_0^T k_i \mu_i X_{\mathbf{k}}^r(u) du \leq c_{11} r^{s_2-p-\varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2-\varepsilon'} \right),$$

for some positive constant c_{11} . For notational convenience, we use \tilde{S} to denote the LHS of the preceding inequality, i.e.,

$$\tilde{S} = \sum_{i:\mathbf{k}+\mathbf{e}_i \in \mathcal{K}} \int_0^T \lambda_i r \frac{X_{\mathbf{k}}^r(u)}{X_{(i)}^r(u)} du + \sum_{i:\mathbf{k}-\mathbf{e}_i \in \bar{\mathcal{K}}} \int_0^T k_i \mu_i X_{\mathbf{k}}^r(u) du, \quad (104)$$

and we have

$$\mathbb{P} \left(\tilde{S} \leq c_{11} r^{s_2-p-\varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2-\varepsilon'} \right) \right) \geq 1 - e^{-\text{poly}(r)}. \quad (105)$$

Let us also recall $S_2(T)$ and $S_3(T)$ defined in (87) and (88), and consider the distribution of $S_2(T) + S_3(T)$. Note that by our construction of the probability space, $S_2(T)$ and $S_3(T)$ are sums of terms that correspond to arrivals and departures driven by independent underlying Poisson processes. Since we are only interested in the distribution of $S_2(T) + S_3(T)$, at this point consider the following alternative construction of the probability space, where all arrivals and departures that appear in the summations of $S_2(T)$ and $S_3(T)$ are driven by a *common* unit-rate Poisson process $\Pi'(\cdot)$. Then \tilde{S} has the same distribution under the original and alternative constructions; and $S_2(T) + S_3(T)$ under the original construction has the same distribution as $\Pi'(\tilde{S})$ under the alternative construction.

We now complete the proof of Claim 2, making use of mainly (97) and (105).

Let $c = \frac{c_1}{64c_{11}}$ and $\bar{c} = \frac{c_1}{4}$. We first consider case 1 of Claim 2, and suppose that $X_{\mathbf{k}}^r(0) \leq 4cr^{s_2}$. For sufficiently large r , $4cr^{s_2} \geq r^{s_2-\varepsilon'}$. Thus, by (105), with probability $1 - e^{-\text{poly}(r)}$,

$$\tilde{S} \leq c_{11} r^{s_2-p-\varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2-\varepsilon'} \right) \leq c_{11} r^{s_2-p-\varepsilon'} \cdot (8cr^{s_2}) \leq \frac{c_1}{8} r^{2s_2-p-\varepsilon'} = \frac{\bar{c}}{2} r^{2s_2-p-\varepsilon'}.$$

Since $\varepsilon' < s_2 - p/2$, $2s_2 - p - \varepsilon' > 0$, and with probability $1 - e^{-\text{poly}(r)}$, we have

$$\Pi'(\tilde{S}) \leq \Pi' \left(\frac{\bar{c}}{2} r^{2s_2-p-\varepsilon'} \right) \leq \bar{c} r^{2s_2-p-\varepsilon'}.$$

Thus, the event

$$\Pi'(\tilde{S}) \leq \frac{c_1}{4} r^{2s_2-p-\varepsilon'},$$

and then also the event

$$S_2(T) + S_3(T) \leq \frac{c_1}{4} r^{2s_2 - p - \varepsilon'},$$

hold with probability $1 - e^{-\text{poly}(r)}$. By (101), we have that with probability $1 - e^{-\text{poly}(r)}$,

$$\begin{aligned} X_{\mathbf{k}}^r(T) - X_{\mathbf{k}}^r(0) &\geq S_1(T) - S_2(T) - S_3(T) \\ &\geq \frac{c_1}{2} r^{s_1 + s_2 - 1 - \varepsilon'} - \bar{c} r^{2s_2 - p - \varepsilon'} \\ &= \bar{c} r^{2s_2 - p - \varepsilon'}, \end{aligned}$$

where the last equality follows from the fact that $s_1 = s_2 + (1 - p)$. This completes the proof of case 1.

Next, consider case 2, and suppose that $X_{\mathbf{k}}^r(0) \geq 2cr^{s_2}$. For sufficiently large r , $X_{\mathbf{k}}^r(0) \geq 2cr^{s_2} \geq r^{s_2 - \varepsilon'}$. Then, with probability $1 - e^{-\text{poly}(r)}$,

$$\tilde{S} \leq c_{11} r^{s_2 - p - \varepsilon'} \left(X_{\mathbf{k}}^r(0) + r^{s_2 - \varepsilon'} \right) \leq 2c_{11} r^{s_2 - p - \varepsilon'} X_{\mathbf{k}}^r(0) \leq \frac{1}{4} X_{\mathbf{k}}^r(0),$$

and with probability $1 - e^{-\text{poly}(r)}$,

$$\Pi'(\tilde{S}) \leq \Pi'\left(\frac{1}{4} X_{\mathbf{k}}^r(0)\right) \leq \frac{1}{2} X_{\mathbf{k}}^r(0).$$

Thus, with probability $1 - e^{-\text{poly}(r)}$, for every $u \in [0, T]$,

$$\begin{aligned} X_{\mathbf{k}}^r(u) - X_{\mathbf{k}}^r(0) &\geq -S_2(u) - S_3(u) \\ &\geq -S_2(T) - S_3(T) \\ &\geq -\frac{1}{2} X_{\mathbf{k}}^r(0), \end{aligned}$$

where we note that the first two inequalities hold with probability 1. This establishes (95), and completes the proof of case 2. This concludes the proof of Claim 2. \square