

MIT Open Access Articles

Flexible Queueing Architectures

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Tsitsiklis, John N. and Kuang Xu. "Flexible Queueing Architectures." *Operations Research* 65,5 (September-October 2017):1398-1413. © 2017 INFORMS.

As Published: <http://dx.doi.org/10.1287/OPRE.2017.1620>

Publisher: Institute for Operations Research and the Management Sciences (INFORMS)

Persistent URL: <https://hdl.handle.net/1721.1/124853>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike



Flexible Queueing Architectures

John N. Tsitsiklis

LIDS, Massachusetts Institute of Technology, Cambridge, MA 02139, jnt@mit.edu

Kuang Xu

Graduate School of Business, Stanford University, Stanford, CA 94305, kuangxu@stanford.edu

We study a multi-server model with n *flexible* servers and n queues, connected through a bipartite graph, where the level of flexibility is captured by an upper bound on the graph's average degree, d_n . Applications in content replication in data centers, skill-based routing in call centers, and flexible supply chains are among our main motivations.

We focus on the scaling regime where the system size n tends to infinity, while the overall traffic intensity stays fixed. We show that a large capacity region and an asymptotically vanishing queueing delay are simultaneously achievable even under limited flexibility ($d_n \ll n$). Our main results demonstrate that, when $d_n \gg \ln n$, a family of expander-graph-based flexibility architectures has a capacity region that is within a constant factor of the maximum possible, while simultaneously ensuring a diminishing queueing delay for *all* arrival rate vectors in the capacity region. Our analysis is centered around a new class of virtual-queue-based scheduling policies that rely on dynamically constructed job-to-server assignments on the connectivity graph. For comparison, we also analyze a natural family of modular architectures, which is simpler but has provably weaker performance. *

Key words: queueing, flexibility, dynamic matching, resource pooling, expander graph, asymptotics

1. Introduction

At the heart of a number of modern queueing networks lies the problem of allocating processing resources (manufacturing plants, web servers, or call-center staff) to meet multiple types of demands that arrive dynamically over time (orders, data queries, or customer inquiries). It is usually the case that a *fully flexible* or *completely resource-pooled* system, where every unit of processing resource is capable of serving all types of demands, delivers the best possible performance. Our inquiry is, however, motivated by the unfortunate reality that such full flexibility is often infeasible due to overwhelming implementation costs (in the case of a data center) or human skill limitations (in the case of a skill-based call center).

What are the key benefits of flexibility and resource pooling in such queueing networks? Can we harness the same benefits even when the degree of flexibility is *limited*, and how should the

* May 2015; revised October 2016. A preliminary version of this paper appeared at Sigmetrics 2013, [30]; the performance of the architectures proposed in the current paper is significantly better than the one in [30]. This research was supported in part by the NSF under grant CMMI-1234062.

network be designed and operated? These are the main questions that we wish to address. While these questions can be approached from a few different angles, we will focus on the metrics of *capacity region* and *expected queueing delay*; the former measures the system's *robustness* against *demand uncertainties*, i.e., when the arrival rates for different demand types are unknown or likely to fluctuate over time, while the latter is a direct reflection of *performance*. Our main message is positive: in the regime where the system size is large, improvements in both the capacity region and delay are *jointly achievable* even under very limited flexibility, given a proper choice of the architecture (interconnection topology) and scheduling policy.

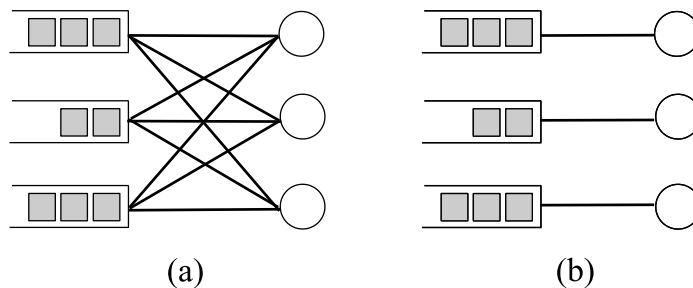


Figure 1 Extreme cases of flexibility: $d_n = n$ versus $d_n = 1$.

Benefits of Full Flexibility. We begin by illustrating the benefits of flexibility and resource pooling in a very simple setting. Consider a system of n servers, each running at rate 1, and n queues, where each queue stores jobs of a particular demand type. For each $i \in \{1, \dots, n\}$, queue i receives an independent Poisson arrival stream of rate λ_i . The average arrival rate $\frac{1}{n} \sum_{i=1}^n \lambda_i$ is denoted by ρ , and is referred to as the *traffic intensity*. The sizes of all jobs are independent and exponentially distributed with mean 1.

For the remainder of this paper, we will use a *measure of flexibility* given by the average number of servers that a demand type can receive service from, denoted by d_n . Let us consider the two extreme cases: a *fully flexible* system, with $d_n = n$ (Figure 1(a)), and an *inflexible* system, with $d_n = 1$ (Figure 1(b)). Fixing the traffic intensity $\rho < 1$, and letting the system size, n , tend to infinity, we observe the following qualitative benefits of full flexibility:

1. Large Capacity Region. In the fully flexible case and under any work-conserving scheduling policy¹, the *collection of all jobs* in the system evolves as an $M/M/n$ queue, with arrival rate $\sum_{i=1}^n \lambda_i$ and service rate n . It is easy to see that the system is stable for all arrival rates that satisfy $\sum_{i=1}^n \lambda_i < n$. In contrast, in the inflexible system, since all $M/M/1$ queues operate independently, we must have $\lambda_i < 1$, for all i , in order to achieve stability. Comparing the two, we see that the

¹ A work-conserving policy mandates that a server be always busy whenever there is at least one job in some queue to which it is connected.

fully flexible system has a much larger capacity region, and is hence more robust to uncertainties or changes in the arrival rates.

2. Diminishing Delay. Let W be the steady-state expected waiting time in queue (time from entering the queue to the initiation of service). As mentioned earlier, the total number jobs in the system for the fully flexible case evolves as an $M/M/n$ queue with traffic intensity $\rho < 1$. It is not difficult to verify that for any fixed value of ρ , the expected total number of jobs in the queues is *bounded above* by a constant independent of n , and hence the expected waiting time in queue satisfies $\mathbb{E}(W) \rightarrow 0$, as $n \rightarrow \infty$.² In contrast, the inflexible system is simply a collection of n independent $M/M/1$ queues, and hence the expected waiting time is $\mathbb{E}(W) = \frac{\rho}{1-\rho} > 0$, for all n . Thus, the expected delay in the fully flexible system *vanishes asymptotically* as the system size increases, but stays bounded away from zero in the inflexible system.

Preview of Main Results. Will the above benefits of fully flexible systems continue to be present if the system only has limited flexibility, that is, if $d_n \ll n$? The main results of this paper show that a large capacity region and an asymptotically vanishing delay can still be *simultaneously achieved*, even when $d_n \ll n$. However, when flexibility is limited, the architecture and scheduling policy need be chosen with care. We show that, when $d_n \gg \ln n$, a family of expander-graph-based flexibility architectures has the largest possible capacity region, up to a constant factor, while simultaneously ensuring a diminishing queueing delay, of order $\ln n/d_n$ as $n \rightarrow \infty$, for *all* arrival rate vectors in the capacity region (Theorem 3.4). For comparison, we also analyze a natural family of modular architectures, which is simpler but has provably weaker performance (Theorems 3.5 and 3.6).

1.1. Motivating Applications

We describe here several motivating applications for our model; Figure 2 illustrates the overall architecture that they share. **Content replication** is commonly used in data centers for bandwidth intensive operations such as database queries [27] or video streaming [20], by hosting the same piece of content on multiple servers. Here, a server corresponds to a physical machine in the data center, and each queue stores incoming demands for a particular piece of content (e.g., a video clip). A server j is connected to queue i if there is a copy of content i on server j , and d_n reflects the average number of replicas per piece of content across the network. Similar structures also arise in **skill-based routing in call centers**, where agents (servers) are assigned to answer calls from different categories (queues) based on their domains of expertise [32], and in **process-flexible**

²The fact that the expected waiting time vanishes asymptotically follows from the bounded expected total number of jobs in steady-state, the assumption that the total arrival rate is ρn , which goes to infinity as $n \rightarrow \infty$, and Little's Law.

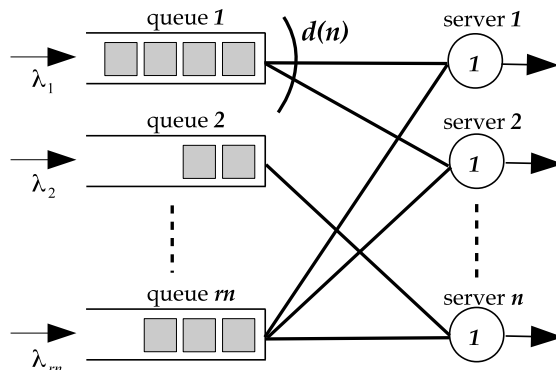


Figure 2 A processing network with rn queues and n servers.

supply chains [16, 26, 6, 15, 10], where each plant (server) is capable of producing multiple product types (queues). In many of these applications, demand rates can be unpredictable and may change significantly over time; for instance, unexpected “spikes” in demand traffic are common in modern data centers [17]. These demand uncertainties make *robustness* an important criterion for system design. These practical concerns have been our primary motivation for studying the interplay between robustness, performance, and the level of flexibility.

1.2. Related Research

Bipartite graphs provide a natural model for capturing the relationships between demand types and service resources. It is well known in the supply chain literature that limited flexibility, corresponding to a sparse bipartite graph, can be surprisingly effective in resource allocation even when compared to a fully flexible system [16, 10, 15, 6, 26]. The use of sparse random graphs or expanders as flexibility structures to improve robustness has recently been studied in [7, 5] in the context of supply chains, and in [20] for content replication. Similar to the robustness results reported in this paper, these works show that random graphs or expanders can accommodate a large set of demand rates. However, in contrast to our work, nearly all analytical results in this literature focus on static allocation problems, where one tries to match supply with demand in a single shot, as opposed to our model, where resource allocation decisions need to be made dynamically over time.

In the queueing theory literature, the models that we consider fall under the umbrella of multi-class multi-server systems, where a set of servers are connected to a set of queues through a bipartite graph. Under these (and similar) settings, complete resource pooling (full flexibility) is known to improve system performance [21, 12, 3]. However, much less is known when only limited flexibility is available: systems with a non-trivial connectivity graph are extremely difficult to analyze, even under seemingly simple scheduling policies (e.g, first-come first-serve) [28, 31]. Simulations in [32] show empirically that limited cross-training can be highly effective in a large call center under

a skill-based routing algorithm. Using a very different set of modeling assumptions, [2] proposes a specific chaining structure with limited flexibility, which is shown to perform well under heavy traffic. Closer to the spirit of the current work is [29], which studies a partially flexible system where a fraction $p > 0$ of all processing resources are fully flexible, while the remaining fraction, $1 - p$, is dedicated to specific demand types, and which shows an exponential improvement in delay scaling under heavy-traffic. However, both [2] and [29] focus on the heavy-traffic regime, which is different from the current setting where traffic intensity is assumed to be fixed, and the analytical results in both works apply only to uniform demand rates. Furthermore, with a constant fraction of the resources being fully flexible, the average degree in [29] scales linearly with the system size n , whereas here we are interested in the case of a much slower (sub-linear) degree scaling.

At a higher level, our work is focused on the interplay between robustness, delay, and the degree of flexibility in a queueing network, which is much less studied in the existing literature, and especially for networks with a non-trivial interconnection topology.

On the technical end, we build on several existing ideas. The techniques of batching (cf. [24, 25]) and the use of virtual queues (cf. [22, 19]) have appeared in many contexts in queueing theory, but the specific models considered in the literature bear little resemblance to ours. The study of expander graphs has become a rich field in mathematics (cf. [14]), but we will refrain from providing a thorough review because only some elementary and standard properties of expander graphs are used in the current paper.

We finally note that preliminary (and weaker) versions of some of the results were included in the conference paper [30].

Organization of the Paper. We describe the model in Section 2, along with the notation to be used throughout. The main results are provided in Section 3. The construction and the analysis associated with the Expander architecture will be presented separately, in Section 4. We conclude the paper in Section 5 with a further discussion of the results as well as directions for future research.

2. Model and Metrics

2.1. Queueing Model and Interconnection Topologies

The Model. We consider a sequence of systems operating in *continuous time*, indexed by the integer n , where the n th system consists of rn queues and n servers (Figure 2), and where r is a constant that is held fixed as n varies. For simplicity, we will set r to 1 but note that all results and arguments in this paper can be extended to the case of general r without difficulty.

A *flexible architecture* is represented by an $n \times n$ undirected bipartite graph $g_n = (E, I \cup J)$, where I and J represent the sets of queues and servers, respectively, and E the set of edges between

them.³ We will also refer to I and J as the sets of left and right nodes, respectively. A server $j \in J$ is *capable* of serving a queue $i \in I$, if and only if $(i, j) \in E$. We will use the following notation.

1. Let \mathcal{G}_n be the set of all $n \times n$ bipartite graphs.
2. For $g_n \in \mathcal{G}_n$, let $\deg(g_n)$ be the average degree among the n left nodes, which is the same as the average degree of the right nodes.
3. For a subset of nodes, $M \subset I \cup J$, let $g|_M$ be the graph induced by g on the nodes in M .
4. Denote by $\mathcal{N}(i)$ the set of servers in J connected to queue i , and similarly, by $\mathcal{N}(j)$ the set of queues in I connected to server j .

Each queue i receives a stream of incoming jobs according to a Poisson process of rate $\lambda_{n,i}$, independent of all other streams, and we define $\boldsymbol{\lambda}_n = (\lambda_{n,1}, \lambda_{n,2}, \dots, \lambda_{n,n})$, which is the **arrival rate vector**. When the value of n is clear from the context, we sometimes suppress the subscript n and write $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$ instead. The sizes of the jobs are exponentially distributed with mean 1, independent from each other and from the arrival processes. All servers are assumed to be running at a constant rate of 1. The system is assumed to be empty at time $t = 0$.

Jobs arriving at queue i can be assigned (immediately, or in the future) to an idle server $j \in \mathcal{N}(i)$ to receive service. The assignment is *binding*: once the assignment is made, the job cannot be transferred to, or simultaneously receive service from, any other server. Moreover, service is *non-preemptive*: once service is initiated for a job, the assigned server has to dedicate its full capacity to this job until its completion.⁴ Formally, if a server j has just completed the service of a previous job at time t or is idle, its available actions are: **(a) Serve a new job**: Server j can choose to fetch a job from any queue in $\mathcal{N}(j)$ and immediately start service. The server will remain occupied and take no other actions until the processing of the current job is completed, which will take an amount of time that is equal to the size of the job. **(b) Remain idle**: Server j can choose to remain idle. While in the idling state, it will be allowed to initiate a service (Action (a)) at any point in time.

Given the limited set of actions available to the server, the performance of the system is fully determined by a *scheduling policy*, π , which specifies for each server $j \in J$, (a) when to remain idle, and when to serve a new job, and (b) from which queue in $\mathcal{N}(j)$ to fetch a job when initiating a new service. We only allow policies that are causal, in the sense that the decision at time t depends only on the history of the system (arrivals and service completions) up to t . We allow the

³ For simplicity of notation, we omit the dependence of E , I , and J on n .

⁴ While we restrict to binding and non-preemptive scheduling policies, other common architectures where (a) a server can serve multiple jobs concurrently (processor sharing), (b) a job can be served by multiple servers concurrently, or (c) job sizes are revealed upon entering the system, are clearly more powerful than the current setting, and are therefore capable of implementing the scheduling policies considered in this paper. As a result, the performance upper bounds developed in this paper also apply to these more powerful variations.

scheduling policy to be *centralized* (i.e., to have full control over all server actions) based on the knowledge of all *queue lengths* and server states. On the other hand, the policy does *not* observe the actual sizes of the jobs before they are served.

2.2. Performance Metrics

Characterization of Arrival Rates. We will restrict ourselves to arrival rate vectors with average *traffic intensity* at most ρ , i.e.,

$$\sum_{i=1}^n \lambda_i \leq \rho n, \quad (1)$$

where $\rho \in (0,1)$ will be treated throughout the paper as a given absolute constant. To quantify the *level of variability* or *uncertainty* of a set of arrival rate vectors, $\mathbf{\Lambda}$, we introduce a *fluctuation parameter*, denoted by u_n , with the property that $\lambda_i < u_n$, for all i and $\boldsymbol{\lambda} \in \mathbf{\Lambda}$.

Note that, for a graph with maximum degree d_n , the fluctuation parameter should not exceed d_n , because otherwise there could exist some $\boldsymbol{\lambda} \in \mathbf{\Lambda}$ under which at least one queue would be unstable. Therefore, the best we can hope for is a flexible architecture that can accommodate arrival rate vectors with a u_n that is close to d_n . The following condition formally characterizes the range of arrival rate vectors we will be interested in, parameterized by the fluctuation parameter, u_n , and traffic intensity, ρ .

Condition 2.1 (Rate Condition) Fix $n \geq 1$ and some $u_n > 0$. We say that a (non-negative) arrival rate vector $\boldsymbol{\lambda}$ satisfies the rate condition if the following hold:

1. $\max_{1 \leq i \leq n} \lambda_i < u_n$.
2. $\sum_{i=1}^n \lambda_i \leq \rho n$.

We denote by $\mathbf{\Lambda}_n(u_n)$ the set of all arrival rate vectors that satisfy the above conditions.

Capacity Region. The capacity region for a given architecture is defined as the set of all arrival rate vectors that it can handle. As mentioned in the Introduction, a larger capacity region indicates that the architecture is more robust against uncertainties or changes in the arrival rates. More formally, we have the following definition.

Definition 2.2 (Feasible Demands and Capacity Region) Let $g = (I \cup J, E)$ be an $n \times n'$ bipartite graph. An arrival rate vector $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$, is said to be *feasible* if there exists a flow, $\mathbf{f} = \{f_{ij} : (i, j) \in E\}$, such that

$$\begin{aligned} \lambda_i &= \sum_{j \in \mathcal{N}(i)} f_{ij}, \quad \forall i \in I, \\ \sum_{i \in \mathcal{N}(j)} f_{ij} &< 1, \quad \forall j \in J, \\ f_{ij} &\geq 0, \quad \forall (i, j) \in E. \end{aligned} \quad (2)$$

In this case, we say that the flow \mathbf{f} satisfies the demand $\boldsymbol{\lambda}$. The capacity region of g , denoted by $\mathbf{R}(g)$, is defined as the set of all feasible demand vectors of g .

It is well known that there exists a policy under which the steady-state expected delay is finite if and only if $\boldsymbol{\lambda} \in \mathbf{R}(g_n)$; the strict inequalities in Definition 2.2 are important here. For the remainder of the paper, we will use the fluctuation parameter u_n (cf. Condition 2.1) to gauge the size of the capacity region, $\mathbf{R}(g_n)$, of an architecture. For instance, if $\boldsymbol{\Lambda}_n(u_n) \subset \mathbf{R}(g_n)$, then the architecture g_n , together with a suitable scheduling policy, allows for finite steady-state expected delay, for any arrival rate vector in $\boldsymbol{\Lambda}_n(u_n)$.

Vanishing Delay. We define the *expected average delay*, $\mathbb{E}(W|\boldsymbol{\lambda}, g, \pi)$ under the arrival rate vector $\boldsymbol{\lambda}$, flexible architecture g , and scheduling policy π , as follows. We denote by $W_{i,m}$ the waiting time in queue experienced by the m th job arriving to queue i , define

$$\mathbb{E}(W_i) = \limsup_{m \rightarrow \infty} \mathbb{E}(W_{i,m}),$$

and let

$$\mathbb{E}(W|\boldsymbol{\lambda}, g, \pi) = \frac{1}{\sum_{i \in I} \lambda_i} \sum_{i \in I} \lambda_i \mathbb{E}(W_i). \quad (3)$$

In the sequel, we will often omit the mention of π , and sometimes of g , and write $\mathbb{E}(W|\boldsymbol{\lambda}, g)$ or $\mathbb{E}(W|\boldsymbol{\lambda})$, in order to place emphasis on the dependencies that we wish to focus on.⁵

The delay performance of the system is measured by the following criteria: (a) for what ranges $\boldsymbol{\Lambda}_n(u_n)$ of arrival rates, $\boldsymbol{\lambda}$, does delay diminish to zero as the system size increases, i.e., $\sup_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}_n(u_n)} \mathbb{E}(W|\boldsymbol{\lambda}) \rightarrow 0$ as $n \rightarrow \infty$, and (b) at what *speed* does the delay diminish, as a function of n ?

2.3. Notation

We will denote by \mathbb{N} , \mathbb{Z}_+ , and \mathbb{R}_+ , the sets of natural numbers, non-negative integers, and non-negative reals, respectively. The following short-hand notation for asymptotic comparisons will be used often, as an alternative to the usual $\mathcal{O}(\cdot)$ notation; here f and g are positive functions, and L is a certain limiting value of interest, in the set of extended reals, $\mathbb{R} \cup \{-\infty, +\infty\}$:

1. $f(x) \lesssim g(x)$ or $g(x) \gtrsim f(x)$ for $\limsup_{x \rightarrow L} f(x)/g(x) < \infty$;
2. $f(x) \ll g(x)$ or $g(x) \gg f(x)$ for $\limsup_{x \rightarrow L} f(x)/g(x) = 0$;
3. $f(x) \sim g(x)$ for $\lim_{x \rightarrow L} f(x)/g(x) = 1$.

We will minimize the use of floors and ceilings, to avoid the cluttering of notation, and thus assume that all values of interest are appropriately rounded up or down to an integer, whenever doing so does not cause ambiguity or confusion. Whenever suitable, we will use upper-case letters for random variables, and lower-case letters for deterministic values.

⁵ Note that $\mathbb{E}(W|\boldsymbol{\lambda})$ captures a *worst-case* expected waiting time across all jobs in the long run, and is always well defined, even under scheduling policies that do not induce a steady-state distribution.

3. Main Results: Capacity Region and Delay of Flexible Architectures

The statements of our main results are given in this section. Below is a high-level summary of our results; a more complete comparison is given in Table 1.

Flexible architectures	Rate Conditions	Capacity Region	Delay
Expander (Theorem 3.4)	$d_n \gg \ln n,$ $u_n \lesssim d_n$	Good for all λ	Good for all λ , with $\mathbb{E}(W) \lesssim \ln n/d_n,$
Modular (Theorems 3.5, 3.7)	$d_n \gg 1,$ $u_n > 1$	Bad for some λ (even if $u_n \lesssim 1$)	Good for uniform λ , with $\mathbb{E}(W) \lesssim \exp(-c \cdot d_n)$
Random Modular (w.h.p.) (Theorems 3.6, 3.7)	$d_n \gtrsim \ln n,$ $u_n \lesssim d_n/\ln n$	Good for most λ , Bad for some λ	Good for most λ , with $\mathbb{E}(W) \lesssim \exp(-c \cdot d_n),$ Bad for some λ

Table 1

This table summarizes and compares the flexibility architectures that we study, in terms of of capacity and delay. We say that capacity is “good” for λ if λ falls within the capacity region of the architecture, and that delay is “good” if the expected delay is vanishingly small for large n . When describing the size of the set of λ for which a statement applies, we use the following (progressively weaker) quantifiers:

1. “For all” means that the statement holds for all $\lambda \in \Lambda_n(u_n)$;
2. “For most” means that the statement holds with high probability when λ is drawn from an arbitrary distribution over $\Lambda_n(u_n)$, independently from any randomization in the construction of the flexibility architecture;
3. “For some” means that the statement is true for a non-empty set of values of λ .

The label “w.h.p.” means that all statements in the corresponding row hold with high probability with respect to the randomness in generating the flexibility architecture.

Our main results focus on an Expander architecture, where the interconnection topology is an expander graph with appropriate expansion. We show that, when $d_n \gg \ln n$, the Expander architecture has a capacity region that is within a constant factor of the maximum possible among all graphs with average degree d_n , while simultaneously ensuring an asymptotically diminishing queueing delay of order $\ln n/d_n$ for all arrival rate vectors in the capacity region, as $n \rightarrow \infty$ (Theorem 3.4). Our analysis involves on a new class of virtual-queue-based scheduling policies that rely on dynamically constructed job-to-server assignments on the connectivity graph.

Our secondary results concern a Modular architecture, which has a simpler construction and scheduling rule compared to the Expander architecture. The Modular architecture consists of a

collection of *separate* smaller subnetworks, with complete connectivity between all queues and servers within each subnetwork. Since the subnetworks are disconnected from each other, a Modular architecture does *not* admit a large capacity region: there always exists an *infeasible* arrival rate vector even when the fluctuation parameter is of constant order (Theorem 3.5). Nevertheless, we show that with proper randomization in the construction of the subnetworks (Randomized Modular architecture), a simple greedy scheduling policy is able to deliver asymptotically vanishing delay for “most” arrival rate vectors with nearly optimal fluctuation parameters, with high probability (Theorem 3.6). These findings suggest that, thanks to its simplicity, the Randomized Modular architecture could be a viable alternative to the Expander architecture if the robustness requirement is not as stringent and one is content with probabilistic guarantees on system stability.

3.1. Preliminaries

Before proceeding, we provide some information on expander graphs, which will be used in some of our constructions and proofs.

Definition 3.1 *An $n \times n'$ bipartite graph $(I \cup J, E)$ is an (α, β) -expander, if for all $S \subset I$ that satisfy $|S| \leq \alpha n$, we have that $|\mathcal{N}(S)| \geq \beta |S|$, where $\mathcal{N}(S) = \bigcup_{i \in S} \mathcal{N}(i)$ is the set of nodes in J that are connected to some node in S .*

The usefulness of expanders in our context comes from the following lemma, which relates the parameters of an expander to the size of its capacity region, as measured by the fluctuation parameter, u_n . The proof is elementary and is given in Appendix A.1.

Lemma 3.2 (Capacity of Expanders) *Fix $n, n' \in \mathbb{N}$, $\rho \in (0, 1)$, $\gamma > \rho$. Suppose that an $n \times n'$ bipartite graph, g_n , is a $(\gamma/\beta_n, \beta_n)$ -expander, where $\beta_n \geq u_n$. Then $\mathbf{A}_n(u_n) \subset \mathbf{R}(g_n)$.*

The following lemma ensures that such expander graphs exist for the range of parameters that we are interested in. The lemma is a simple consequence of a standard result on the existence of expander graphs, and its proof is given in Appendix A.2.

Lemma 3.3 *Fix $\rho \in (0, 1)$. Suppose that $d_n \rightarrow \infty$ as $n \rightarrow \infty$. Let $\beta_n = \frac{1}{2} \cdot \frac{\ln(1/\rho)}{1 + \ln(1/\rho)} d_n$, and $\gamma = \sqrt{\rho}$. There exists $n' > 0$, such that for all $n \geq n'$, there exists an $n \times n$ bipartite graph which is a $(\gamma/\beta_n, \beta_n)$ -expander with maximum degree d_n .*

Remark. It is well known that random graphs with appropriate average degree are expanders with high probability (cf. [14]). For instance, it is not difficult to show that if $d_n \gg \ln n$ and $\beta_n = \frac{1-\gamma}{4} d_n / \ln n$, then an Erdős-Rényi random bipartite graph with average degree d_n is a $(\gamma/\beta_n, \beta_n)$ -expander, with high probability, as $n \rightarrow \infty$ (cf. Lemma 3.12 of [33]). We note, however, that to

deterministically construct expanders in a computationally efficient manner can be challenging and is in and of itself an active field of research; the reader is referred to the survey paper [14] and the references therein.

3.2. Expander Architecture

Construction of the Architecture. The connectivity graph in the Expander Architecture is an expander graph with maximum degree d_n and appropriate expansion.

Scheduling Policy. We employ a scheduling policy that organizes the arrivals into batches, stores the batches in a virtual queue, and dynamically assigns the jobs in a batch to appropriate servers. Theorem 3.4, which is the main result of this paper, shows that under this policy the Expander architecture achieves an asymptotically *vanishing* delay for *all* arrival rate vectors in the set $\Lambda_n(u_n)$. Of course we assume that d_n is sufficiently large so that the corresponding expander graph exists (Lemma 3.3, with ρ replaced with $\hat{\rho}$). At a high level, the strong guarantees stem from the excellent connectivity of an expander graph, and similarly of random subsets of an expander graph, a fact which we will exploit to show that jobs arriving to the system during a small time interval can be quickly assigned to connected idle servers with high probability, which then leads to a small delay. The proof of the theorem, including a detailed description of the scheduling policy, is given in Section 4.

Theorem 3.4 (Capacity and Delay of Expander Architectures) *Let $\hat{\rho} = \frac{1}{1+(1-\rho)/8}$. For every $n \in \mathbb{N}$, define*

$$\beta_n = \frac{1}{2} \cdot \frac{\ln(1/\hat{\rho})}{\ln(1/\hat{\rho}) + 1} d_n,$$

and

$$\gamma = \sqrt{\hat{\rho}}.$$

Suppose that $\ln n \ll d_n \ll n$, and

$$u_n \leq \frac{1-\rho}{2} \beta_n.$$

Let g_n be a $(\gamma/\beta_n, \beta_n)$ -expander with maximum degree d_n . The following holds.

1. There exists a scheduling policy, π_n , under which

$$\sup_{\lambda_n \in \Lambda_n(u_n)} \mathbb{E}(W | \lambda_n, g_n) \leq \frac{c \ln n}{d_n}, \quad (4)$$

where c is a constant independent of n and g_n .

2. The scheduling policy, π_n , only depends on g_n and an upper bound on the traffic intensity, ρ . It does not require knowledge of the arrival rate vector λ_n .

Note that when ρ is viewed as a constant, the upper bound on u_n in the statement of Theorem 3.4 is just a constant multiple of d_n . Since the fluctuation parameter, u_n , should be no more than d_n for stability to be possible, the size of $\Lambda_n(u_n)$ in Theorem 3.4 is within a constant factor of the *best possible*.

Remark. Compared to our earlier results, in a preliminary version of this paper (Theorem 1 in [30]), Theorem 3.4 is stronger in two major aspects: (1) the guarantee for diminishing delay holds deterministically over all arrival rate vectors in $\Lambda_n(u_n)$, as opposed to “with high probability” over the randomness in the generation of g_n , and (2) the fluctuation parameter, u_n , is allowed to be of order d_n in Theorem 3.4, while [30] required that $u_n \ll \sqrt{d_n/\ln n}$. The flexible architecture considered in [30] was based on Erdős-Rényi random graphs. It also employed a scheduling policy based on virtual queues, as in this paper. However, the policy in the present paper is simpler to describe and analyze.

3.3. Modular Architectures

In a Modular architecture, the designer partitions the network into n/d_n *separate* subnetworks. Each subnetwork consists of d_n queues and servers that are fully connected (Figure 3), but disconnected from queues and servers in other subnetworks.

Construction of the Architecture. Formally, the construction is as follows.

1. We partition the set J of servers into n/d_n disjoint subsets (“clusters”) $B_1, \dots, B_{n/d_n}$, all having the same cardinality d_n . For concreteness, we assign the first d_n servers to the first cluster, B_1 , the next d_n servers to the second cluster, etc.
2. We form a partition $\sigma_n = (A_1, \dots, A_{n/d_n})$ of the set I of queues into n/d_n disjoint subsets (“clusters”) A_k , all having the same cardinality d_n .
3. To construct the interconnection topology, for $k = 1, \dots, n/d_n$, we connect every queue $i \in A_k$ to every server $j \in B_k$. A pair of queue and server clusters with the same index k will be referred to as a subnetwork.

Note that in a Modular architecture, the degree of each node is equal to the size, d_n , of the clusters. Note also that different choices of σ_n yield isomorphic architectures. When σ_n is drawn uniformly at random from the set of all possible partitions of I into subsets of size n/d_n , we call the resulting topology a *Random Modular* architecture.

Scheduling Policy. We use a simple greedy policy, equivalent to running each subnetwork as an $M/M/d_n$ queue. Whenever a server $j \in B_k$ becomes available, it starts serving a job from any non-empty queue in A_k . Similarly, when a job arrives at queue $i \in A_k$, it is immediately assigned to an arbitrary idle server in B_k , if such a server exists, and waits in queue i , otherwise.

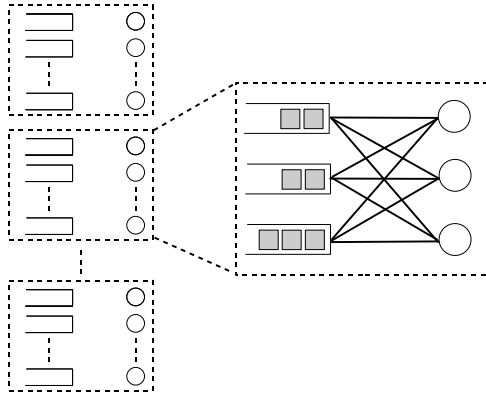


Figure 3 A Modular architecture consisting of n/d_n subnetworks, each with d_n queues and servers. Within each subnetwork, all servers are connected to all queues.

Our first result points out that a Modular architecture does not have a large capacity region: for any partition σ_n , there always exists an infeasible arrival rate vector, even if u_n is small, of order $\mathcal{O}(1)$. The proof is given in Appendix A.3. Note that this is a negative result that applies no matter what scheduling policy is used.

Theorem 3.5 (Capacity Region of Deterministic Modular Architectures) *Fix $n \geq 1$ and some $u_n > 1$. Let g_n be a Modular architecture with average degree $d_n \leq \frac{\rho}{2}n$. Then, there exists $\lambda_n \in \Lambda_n(u_n)$ such that $\lambda_n \notin \mathbf{R}(g_n)$.*

However, if we are willing to settle for a weaker result on the capacity region, the next theorem states that with the Random Modular architecture, any given arrival rate vector λ_n has high probability (with respect to the random choice of the partition σ_n) of belonging to the capacity region, if the fluctuation parameter, u_n , is of order $\mathcal{O}(d_n/\ln n)$, but no more than that. Intuitively, this is because the randomization in the connectivity structure makes it unlikely that many large components of λ_n reside in the same sub-network. The proof is given in Appendix A.4.

Theorem 3.6 (Capacity Region of Random Modular Architectures) *Let σ_n be drawn uniformly at random from the set of all partitions, and let G_n be the resulting Random Modular architecture. Let \mathbb{P}_{G_n} be the probability measure that describes the distribution of G_n . Fix a constant $c_1 > 0$, and suppose that $d_n \geq c_1 \ln n$. Then, there exist positive constants c_2 and c_3 , such that:*

(a) *If $u_n \leq c_2 d_n / \ln n$, then*

$$\lim_{n \rightarrow \infty} \inf_{\lambda_n \in \Lambda_n(u_n)} \mathbb{P}_{G_n}(\lambda_n \in \mathbf{R}(G_n)) = 1. \quad (5)$$

(b) *Conversely, if $u_n > c_3 d_n / \ln n$ and $d_n \leq n^{0.3}$, then*

$$\lim_{n \rightarrow \infty} \inf_{\lambda_n \in \Lambda_n(u_n)} \mathbb{P}_{G_n}(\lambda_n \in \mathbf{R}(G_n)) = 0, \quad (6)$$

We can use Theorem 3.6 to obtain a statement about “most” arrival rate vectors in $\Lambda_n(u_n)$, as follows. Suppose that λ_n is drawn from an arbitrary distribution μ_n over $\Lambda_n(u_n)$, independently from the randomness in G_n . Let $\mathbb{P}_{G_n} \times \mu_n$ be the product measure that describes the joint distribution of G_n and λ_n . Using Fubini’s theorem, Eq. (5) implies that

$$\lim_{n \rightarrow \infty} (\mathbb{P}_{G_n} \times \mu_n)(\lambda_n \in \mathbf{R}(G_n)) = 1. \quad (7)$$

A further application of Fubini’s theorem and an elementary argument⁶ implies that there exists a sequence δ'_n that converges to zero, such that the event

$$\mu_n(\lambda_n \in \mathbf{R}(G_n) \mid G_n) \geq 1 - \delta'_n \quad (8)$$

has “high probability,” with respect to the measure \mathbb{P}_{G_n} . That is, there is high probability that the Random Modular architecture includes “most” arrival vectors λ_n in $\Lambda_n(u_n)$.

We now turn to delay. The next theorem states that in a Modular architecture, delay is vanishingly small for all arrival rate vectors in the capacity region that are not too close to its outer boundary. The proof is given in Appendix A.5.

We need some notation. For any set S and scalar γ , we let $\gamma S = \{\gamma x : x \in S\}$.

Theorem 3.7 (Delay of Modular Architectures) *Fix some $\gamma \in (0, 1)$, and consider a Modular architecture g_n for each n . There exists a constant $c > 0$, independent of n and the sequence $\{g_n\}$, so that*

$$\mathbb{E}(W \mid \lambda_n) \lesssim \exp(-c \cdot d_n), \quad (9)$$

for every $\lambda_n \in \gamma \mathbf{R}(g_n)$.

3.3.1. Expanded Modular Architectures There is a further variant of the Modular architecture that we call the *Expanded Modular architecture*, which combines the features of a Modular architecture and an expander graph via a graph product. By construction, it uses part of the system flexibility to achieve a large capacity region and part to achieve low delay. As a result, the Expanded Modular architecture admits a smaller capacity region compared to that of an Expander architecture. Another drawback is that the available performance guarantees involve policies that require the knowledge of the arrival rates λ_i . On the positive side, it guarantees an asymptotically vanishing delay for *all* arrival rates, uniformly across the capacity region, and can be operated by a scheduling policy that is arguably simpler than in the Expander architecture. The construction and a scheduling policy for the Expanded Modular architecture is given in Appendix B, along with a statement of its performance guarantees (Theorem B.1). The technical details can be found in [33].

⁶ We are using here the following elementary Lemma. Let A be an event with $\mathbb{P}(A) \geq 1 - \epsilon$, and let X be a random variable. Then, there exists a set B with $\mathbb{P}(B) \geq 1 - \sqrt{\epsilon}$ such that $\mathbb{P}(A \mid X) \geq 1 - \sqrt{\epsilon}$, whenever $X \in B$. The lemma is applied by letting A be the event $\{\lambda_n \in \mathbf{R}(G_n)\}$ and letting $X = G_n$.

4. Analysis of the Expander Architecture

In this section, we introduce a policy for the Expander architecture, based on batching and virtual queues, which will then be used to prove Theorem 3.4. We begin by describing the basic idea at a high level.

4.1. The Main Idea

Our policy proceeds by collecting a fair number of arriving jobs to form batches. Batches are thought of as being stored in a virtual queue, with each batch treated as a single entity. By choosing the batch size large enough, one expects to see certain statistical regularities that can be exploited in order to efficiently handle the jobs within a batch. We now provide an outline of the operation of the policy, for a special case.

Let us fix n and consider the case where $\lambda_i = \lambda < 1$ for all i . Suppose that at time t , all servers are busy serving some job. Let us also fix some γ_n such that $\gamma_n \ll 1$, while $n\gamma_n$ is large. During the time interval $[t, t + \gamma_n)$, “roughly” $\lambda n\gamma_n$ new jobs will arrive and $n\gamma_n$ servers will become available. Let Γ be the set of queues that received any job and let Δ be the set of servers that became available during this interval. Since $\lambda n\gamma_n \ll n$, these incoming jobs are likely to be spread out across different queues, so that most queues receive at most one job. Assuming that this is indeed the case, we focus on $g_n|_{\Gamma \cup \Delta}$, that is, the connectivity graph g_n , restricted to $\Gamma \cup \Delta$. The key observation is that this is a subgraph sampled uniformly at random among all subgraphs of g_n with approximately $\lambda n\gamma_n$ left nodes and $n\gamma_n$ right nodes. When $n\gamma_n$ is sufficiently large, and g_n is well connected (as in an expander with appropriate expansion properties), we expect that, with high probability, $g_n|_{\Gamma \cup \Delta}$ admits a matching that includes the entire set Γ (i.e., a one-to-one mapping from Γ to Δ). In this case, we can ensure that *all* of the roughly $\lambda n\gamma_n$ jobs can start receiving service at the end of the interval, by assigning them to the available servers in Δ according to this particular matching. Note that the resulting queueing delay will be comparable to γ_n , which has been assumed to be small.

The above described scenario corresponds to the normal course of events. However, with a small probability, the above scenario may not materialize, due to statistical fluctuations, such as:

1. Arrivals may be concentrated on a small number of queues.
2. The servers that become available may be located in a subset of g_n that is not well connected to the queues with arrivals.

In such cases, it may be impossible to assign the jobs in Γ to servers in Δ . These exceptional cases will be handled by the policy in a different manner. However, if we can guarantee that the probability of such cases is low, we can then argue that their impact on performance is negligible.

Whether or not the above mentioned exceptions will have low probability of occurring depends on whether the underlying connectivity graph, g_n , has the following property: with high probability, a randomly sampled sublinear (but still sufficiently large) subgraph of g_n admits a large set of “flows.” This property will be used to guarantee that, with high probability, the jobs in Γ can indeed be assigned to distinct servers in the set Δ . We will show that an expander graph with appropriate expansion does possess this property.

4.2. An additional assumption

Before proceeding, we introduce an additional assumption on the arrival rates, which will remain in effect throughout this section, and which will simplify some of the arguments. Appendix A.6 explains why this assumption can be made without loss of generality.

Assumption 4.1 (Lower Bound on the Total Arrival Rate) *We have that $\rho \in (1/2, 1)$, and the total arrival rate satisfies the lower bound*

$$\sum_{i=1}^n \lambda_i \geq (1 - \rho)n. \quad (10)$$

4.3. The Policy

We now describe in detail the scheduling policy. Besides n , the scheduling policy uses the following inputs:

1. ρ , the traffic intensity introduced in Condition 2.1, in Section 2.2,
2. ϵ , a positive constant such that $\rho + \epsilon < 1$.
3. b_n , a batch size parameter,
4. g_n , the connectivity graph.

Notice that the arrival rates, λ_i , and the fluctuation parameter, u_n , are *not* inputs to the scheduling policy.

At this point it is useful to make a clarification regarding the \lesssim notation. Recall that the relation $f(n) \lesssim g(n)$ means that $f(n) \leq cg(n)$, for all n , where c is a positive constant. Whenever we use this notation, we require that the constant c cannot depend on any parameters other than ρ and ϵ . Because we view ρ and ϵ as fixed throughout, this makes c an absolute constant.

4.3.1. Arrivals of Batches. Arriving jobs are organized in *batches* of cardinality ρb_n , where b_n is a design parameter, to be specified later.⁷ Let $T_0^B = 0$. For $k \geq 1$, let T_k^B be the time of the $(k\rho b_n)$ th arrival to the system, which we also view as the *arrival time* of the k th batch. For $k \geq 1$,

⁷ In a slight departure from the earlier informal description, we define batches by keeping track of the number of arriving jobs as opposed to keeping track of time.

the k th batch consists of the ρb_n jobs that arrive during the time interval $(T_{k-1}^B, T_k^B]$. The length $A_k = T_k^B - T_{k-1}^B$ of this interval will be called the k th *inter-arrival time*. We record, in the next lemma, some immediate statistical properties of the batch inter-arrival times.

Lemma 4.2 *The batch inter-arrival times, $\{A_k\}_{k \geq 1}$, are i.i.d., with*

$$\frac{b_n}{n} \leq \mathbb{E}(A_k) \leq \frac{\rho}{1-\rho} \cdot \frac{b_n}{n},$$

and $\text{Var}(A_k) \lesssim b_n/n^2$.

Proof. The batch inter-arrival times are i.i.d., due to our independence assumptions on the job arrivals. By definition, A_k is equal in distribution to the time until a Poisson process records ρb_n arrivals. This Poisson process has rate $r = \sum_{i=1}^n \lambda_i$, and using also Assumption 4.1 in the first inequality below, we have

$$(1-\rho)n \leq \sum_{i=1}^n \lambda_i = r \leq \rho n,$$

The random variables A_k are Erlang (sum of ρb_n exponentials with rate r). Therefore,

$$\mathbb{E}(A_k) = \rho b_n \cdot \frac{1}{r} \geq \rho b_n \cdot \frac{1}{\rho n} = \frac{b_n}{n}.$$

Similarly,

$$\mathbb{E}(A_k) = \rho b_n \cdot \frac{1}{r} \leq \rho b_n \cdot \frac{1}{(1-\rho)n}.$$

Finally,

$$\text{Var}(A_k) = \rho b_n \cdot \frac{1}{r^2} \leq \rho b_n \cdot \frac{1}{(1-\rho)^2 n^2} \lesssim \frac{b_n}{n^2}.$$

Q.E.D.

4.3.2. The Virtual Queue Upon arrival, batches are placed in what we refer to as a *virtual queue*. The virtual queue is a GI/G/1 queue, which is operated in FIFO fashion. That is, a batch waits in queue until all previous batches are served, and then starts being served by a virtual queueing system. The service of a batch by the virtual queueing system lasts until a certain time by which all jobs in the batch have already been assigned to, and have started receiving service from, one of the physical servers, at which point the service of the batch is completed and the batch departs from the virtual queue. The time elapsed from the initiation of service of batch until its departure is called the *service time* of the batch. As a consequence, the queueing delay of a job in the actual (physical) system is bounded above by the sum of:

- (a) the time from the arrival of the job until the arrival time of the batch that the job belongs to;
- (b) the time that the batch waits in the virtual queue;
- (c) the service time of the batch.

Service slots. The service of the batches at the virtual queue is organized along consecutive time intervals that we refer to as *service slots*. The service slots are intervals of the form $(ls, (l+1)s]$, where l is a nonnegative integer, whose length is⁸

$$s = (\rho + \epsilon) \cdot \frac{b_n}{n}.$$

We will arrange matters so that batches can complete service and depart *only* at the end of a service slot, that is, at times of the form ls . Furthermore, we assume that the physical servers are operated as follows. If either a batch completes service at time ls or if there are no batches present at the virtual queue at that time, we assign to every idle server a dummy job whose duration is an independent exponential random variable, with mean 1. This ensures that the state of the n servers is the same (all of them are busy) at certain special times, thus facilitating further analysis, albeit at the cost of some inefficiency.

4.3.3. The Service Time of a Batch. The specification of the service time of a batch depends on whether the batch, upon arrival, finds an empty or nonempty virtual queue.

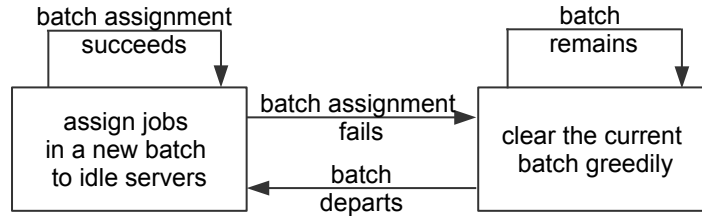


Figure 4 An illustration of the service slot dynamics. An arrow indicates the transition from the end of one service slot to the next.

Suppose that a batch arrives during the service slot $(ls, (l+1)s]$ and finds an empty virtual queue; that is, all previous batches have departed by time ls . According to what was mentioned earlier, at time ls , all physical servers are busy, serving either real or dummy jobs. Up until the end of the service slot, any server that completes service is not assigned a new (real or dummy) job, and remains idle, available to be assigned a job at the very end of the service slot. Let Δ be the set of servers that are idle at time $(l+1)s$, the end of the service slot. At that time, we focus on the jobs in the batch under consideration. We wish to assign each job i in this batch to a *distinct* server $j \in \Delta$, subject to the constraint that $(i, j) \in E$. We shall refer to such a job-to-server assignment as a *batch assignment*. There are two possibilities (cf. Figure 4):

⁸ To see how the length of the service slot was chosen, recall that the size of each batch is equal to ρb_n . The length of the service slot hence ensures that $(\rho + \epsilon)b_n$, the expected number of servers that will become available (and can therefore be assigned to jobs) during a single service slot, is greater than the size of a batch, so that there is hope of assigning all of these jobs to available servers within a single service slot. At the same time, since $\rho + \epsilon < 1$, service slots are shorter than the expected batch inter-arrival time, which is needed for the stability of the virtual queue.

(a) If a batch assignment can be found, each job in the batch is assigned to a server according to that assignment, and the batch departs at time $(l + 1)s$. In this case, we say that the service time of the batch was *short*.

(b) If a batch assignment cannot be found, we start assigning the jobs in the batch to physical servers in some arbitrary greedy manner: Whenever a server j becomes available, we assign to it a job from the batch under consideration, and from some queue i with $(i, j) \in E$, as long as such a job exists. (Ties are broken arbitrarily.) As long as every queue is connected to at least one server, all jobs in the associated batch will be eventually assigned. The last of the jobs in the batch gets assigned during a subsequent service interval $(l's, (l' + 1)s]$, where $l' > l$, and we define $(l' + 1)s$ as the departure time of the batch.

If the k th batch did indeed find an empty virtual queue upon arrival, its service time, denoted by S_k , is the time elapsed from its arrival until its departure.

Suppose now that a batch arrives during a service slot $(ls, (l + 1)s]$ and finds a non-empty virtual queue; that is, there are one or more batches that arrived earlier and which have not departed by time ls . In this case, the batch waits in the virtual queue until some time of the form $l's$, with $l' > l$, when the last of the previous batches departs. Recall that, as specified earlier, at time $l's$ all servers are made to be busy (perhaps, by giving them dummy jobs) and we are faced with a situation identical to the one considered in the previous case, as if the batch under consideration just arrived at time $l's$; in particular, the same service policy can be applied. For this case, where the k th batch arrives to find a non-empty virtual queue, its service time, S_k , extends from the time of the departure of the $(k - 1)$ st batch until the departure of the k th batch.

4.4. Bounding the Virtual Queue by a GI/GI/1 Queue

Having defined the inter-arrival and service times of the batches, the virtual queue is a fully specified, work-conserving, FIFO single-server queueing system.

We note however one complication. The service times of the different batches are dependent on the arrival times. To see this, suppose, for example, that a batch upon arrival sees an empty virtual queue and that its service time is “short.” Then, its service time will be equal to the remaining time until the end of the current service slot, and therefore dependent on the batch’s arrival time. Furthermore, the service times of different batches are dependent: if the service time of the previous batch happens to be too long, then the next batch is likely to see upon arrival a non-empty virtual queue, which then implies that its own service time will be an integer multiple of s .

In order to get around these complications, and to be able to use results on GI/GI/1 queues, we define the *modified service time*, S'_k , of the k th service batch to be equal to S_k , rounded above to the nearest integer multiple of s :

$$S'_k = \min\{ls : ls \geq S_k, l = 1, 2, \dots\}.$$

Clearly, we have $S_k \leq S'_k$.

We now consider a modified (but again FIFO and work-conserving) virtual queueing system in which the arrival times are the same as before, but the service times are the S'_k . A simple coupling argument, based on Lindley's recursion, shows that for every sample path, the time that the batch spends waiting in the queue of the original virtual queueing system is less than or equal to the time spent waiting in the queue of the modified virtual queueing system. It therefore suffices to upper bound the expected time spent in the queue of the modified virtual queueing system.

We now argue that the modified virtual queueing system is a GI/GI/1 queue, i.e., that the service times S'_k are i.i.d., and independent from the arrival process. For a batch whose service starts during the service slot $[ls, (l+1)s)$, the modified service time is equal to s , whenever the batch service time is short. Whether the batch service time will be short or not is determined by the composition of the jobs in this batch and by the identities of the servers who complete service during the service slot $[ls, (l+1)s)$. Because the servers start at the same "state" (all busy) at each service slot, it follows that the events that determine whether a batch service time will be short or not are independent across batches, and with the same associated probabilities.

Similarly, if a batch service time is not short, the additional time to serve the jobs in the batch is affected only by the composition of jobs in the batch and the service completions at the physical servers after time ls , and these are again independent from the inter-arrival times and the modified service times S'_m of other batches m . Finally, the same considerations show the independence of the S'_k from the batch arrival process.

It should now be clear from the above discussion that the modified service time of a batch is of the form

$$S'_k = s + X_k \cdot \hat{S}_k, \quad (11)$$

where:

(a) X_k is a Bernoulli random variable which is equal to 1 if and only if the k th batch service time is not short, i.e., it takes more than a single service slot;

(b) \hat{S}_k is a random variable which (assuming that every queue is connected to at least one server) is stochastically dominated by the sum of ρb_n independent exponential random variables with mean 1, rounded up to the nearest multiple of s . (This dominating random variable corresponds to the extreme case where all of the ρb_n jobs in the batch are to be served in sequence, by the same physical server.)

(c) The pairs (X_k, \hat{S}_k) are i.i.d.

4.5. Bounds on the Modified Service Times

For the remainder of Section 4, we will assume that

$$g_n \text{ is a } (\gamma/\beta_n, \beta_n)\text{-expander,} \quad (12)$$

where γ and β_n are defined as in the statement of Theorem 3.4.

The main idea behind the rest of the proof is as follows. We will upper bound the expected time spent in the modified virtual queueing system using Kingman's bound [18] for GI/GI/1 queues. Indeed, the combination of a batching policy with Kingman's bound is a fairly standard technique for deriving delay upper bounds (see, e.g., [25]). We already have bounds on the mean and variance of the inter-arrival times. In order to apply Kingman's bound, it remains to obtain bounds on the mean and variance of the service times S'_k of the modified virtual queueing system.

We now introduce an important quantity associated with a graph g_n , by defining

$$q(g_n) = \mathbb{P}(X_k = 1 \mid g_n);$$

because of the i.i.d. properties of the batch service times, this quantity does not depend on k . In words, for a given connectivity graph g_n , the quantity $q(g_n)$ stands for the probability that we cannot find a batch assignment, between the jobs in a batch and the servers that become idle during a period of length s .

We begin with the following lemma, which provides bounds on the mean and variance of S'_k .

Lemma 4.3 *There exists a sequence, $\{c_n\}_{n \in \mathbb{N}}$, with $c_n \lesssim b_n$, such that for all $n \geq 1$*

$$s \leq \mathbb{E}(S'_k \mid g_n) \leq s + q(g_n)c_n,$$

and

$$\text{Var}(S'_k \mid g_n) \lesssim q(g_n)c_n^2.$$

Proof. The fact that $\mathbb{E}(S'_k \mid g_n) \geq s$ follows from the definition of S'_k in Eq. (11) and the non-negativity of $X_k \hat{S}_k$. The definition of an expander ensures that every queue is connected to at least one server through g_n . Recall that \hat{S}_k is zero if $X_k = 0$; on the other hand, if $X_k = 1$, and as long as every queue is connected to some server, then \hat{S}_k is upper bounded by the sum of ρb_n exponential random variables with mean 1, rounded up to an integer multiple of s . Therefore,

$$\mathbb{E}(S'_k \mid g_n) = s + \mathbb{E}(X_k \hat{S}_k \mid g_n) = s + \mathbb{P}(X_k = 1 \mid g_n) \cdot \mathbb{E}(\hat{S}_k \mid X_k = 1, g_n) \leq s + q(g_n)(\rho b_n + s),$$

which leads to the first bound in the statement of the lemma, with $c_n = b_n + s$. Since s is proportional to b_n/n , we also have $c_n \lesssim b_n$, as claimed. Furthermore,

$$\text{Var}(S'_k \mid g_n) = \text{Var}(X_k \hat{S}_k \mid g_n) \leq \mathbb{E}(X_k^2 \hat{S}_k^2 \mid g_n) \lesssim q(g_n)(b_n + s)^2 = q(g_n)c_n^2.$$

Q.E.D.

We now need to obtain bounds on $q(g_n)$. This is nontrivial and forms the core of the proof of the theorem. In what follows, we will show that with appropriate assumptions on the various parameters, and for any $\lambda \in \Lambda(u_n)$, an Erdős-Rényi random graph has a very small $q(g_n)$, with high probability.

4.6. Assumptions on the Various Parameters

From now on, we focus on a specific batch size parameter of the form

$$b_n = \frac{320}{(1-\rho)^2} \cdot \frac{n \ln n}{\beta_n}. \quad (13)$$

We shall also set

$$\epsilon = \frac{1-\rho}{2}. \quad (14)$$

We assume, as in the statement of Theorem 3.4, that $d_n \ll n$, and that

$$\beta_n \gtrsim d_n \gg \ln n. \quad (15)$$

Under these choices of b_n and d_n , we have

$$b_n \lesssim \frac{n}{d_n / \ln n} \ll n; \quad (16)$$

that is, the batch size is *vanishingly small* compared to n . Finally, we will only consider arrival rate vectors that belong to the set $\Lambda_n(u_n)$ (cf. Condition 2.1), where, as in the statement of Theorem 3.4,

$$u_n \leq \frac{1-\rho}{2} \beta_n. \quad (17)$$

4.7. The Probability of a Short Batch Service Time

We now come to the core of the proof, aiming to show that if the connectivity graph g_n is an expander graph with a sufficiently large expansion factor, then $q(g_n)$ is small. More precisely, we aim to show that a typical batch will have high probability of having a short service time. A concrete statement is given in the result that follows, and the rest of this subsection will be devoted to its proof.

Proposition 4.4 *Fix $n \geq 1$. We have that*

$$q(g_n) \leq \frac{1}{n^2}. \quad (18)$$

Let us focus on a particular batch, and let us examine what it takes for its service time to be short. There are two sources of randomness:

1. A total of ρb_n jobs arrive to the queues. Let A_i be the number of jobs that arrive at the i th queue, let $\mathbf{A} = (A_1, \dots, A_n)$, and let Γ be the set of queues that receive at least one job. In particular, we have

$$\sum_{i=1}^n A_i = \sum_{i \in \Gamma} A_i = \rho b_n.$$

2. During the time slot at which the service of the batch starts, each server starts busy (with a real or dummy job). With some probability, and independently from other servers or from the arrival process, a server becomes idle by the end of the service time slot. Let Δ be the set of servers that become idle.

Recalling the definition of X_k as the indicator random variable of the event that the service time of the k th batch is not short, we see that X_k is completely determined by the graph g_n together with \mathbf{A} and Δ . For the remainder of this subsection, we suppress the subscript k , since we are focusing on a particular batch. We therefore have a dependence of the form

$$X = f(g_n, \mathbf{A}, \Delta),$$

for some function f , and we emphasize the fact that \mathbf{A} and Δ are independent.

Recall that $\epsilon = (1 - \rho)/2$, and from the statement of Theorem 3.4 that

$$\hat{\rho} = \frac{1}{1 + (1 - \rho)/8} = \frac{1}{1 + \epsilon/4}. \quad (19)$$

Clearly, $\hat{\rho} < 1$, and with some elementary algebra, it is not difficult to show that, for any given $\rho \in (0, 1)$,

$$\hat{\rho} > \rho.$$

Let

$$m_n = \frac{\rho}{\hat{\rho}} b_n,$$

so that

$$\hat{\rho} m_n = \rho b_n. \quad (20)$$

Finally, let

$$\hat{u}_n = \beta_n \frac{m_n}{n}.$$

We will say that \mathbf{A} is *nice* if there exists a set $\hat{\Gamma} \supset \Gamma$ of cardinality m_n , such that $A_i = 0$ whenever $i \notin \hat{\Gamma}$, and

$$A_i < \hat{u}_n, \quad \forall i \in \hat{\Gamma}.$$

We now establish that \mathbf{A} is nice, with high probability. The main idea is simple: \mathbf{A} is not nice only if one out of a finite collection of binomial variables with large means takes a value which is

away from its mean by a certain multiplicative factor. Using the Chernoff bound, this probability can be shown to decay at least as fast as $1/n^3$. The details of this argument are given in the proof of Lemma 4.5, in Appendix A.7.

Lemma 4.5 *For all sufficiently large n , we have that*

$$\mathbb{P}(\mathbf{A} \text{ is not nice}) \leq \frac{1}{n^3}.$$

We now wish to establish that when \mathbf{A} is nice, there is high probability (with respect to Δ), that the batch service time will be short. Having a short batch service time is, by definition, equivalent to the existence of a batch assignment, which in turn is equivalent to the existence of a certain flow in a subgraph of g_n . The lemma that follows deals with the latter existence problem for the original graph, but will be later applied to subgraphs. Let $\overline{\mathbf{R}}(g)$ be the closure of the capacity region, $\mathbf{R}(g)$, of g .

Lemma 4.6 *Fix $n, n' \in \mathbb{N}$, $\rho \in (0, 1)$, and $\gamma > \rho$. Suppose that an $n \times n'$ bipartite graph, g_n , is a $(\gamma/\beta_n, \beta_n)$ -expander, where $\beta_n \geq u_n$. Then $\mathbf{A}_n(u_n) \subset \overline{\mathbf{R}}(g_n)$.*

Proof. The claim follows directly from Lemma 3.2, by noting that $\overline{\mathbf{R}}(g_n) \supset \mathbf{R}(g_n)$. Q.E.D.

The next lemma is the key technical result of this subsection. It states that if g_n is an expander, then, for any given $\hat{\Gamma}$, the random subgraph $g_n|_{\hat{\Gamma} \cup \Delta}$ will be an expander graph with high probability (with respect to Δ). The lemma is stated as a stand-alone result, though we will use a notation that is consistent with the rest of the section. The proof relies on a delicate application of the Chernoff bound, and is given in Appendix A.8.

Lemma 4.7 *Fix $n \geq 1$, $\gamma \in (0, 1)$, and $\rho \in [1/2, 1)$. Let $g_n = (I \cup J, E)$ be an $n \times n$ bipartite graph that is a $(\gamma/\beta_n, \beta_n)$ -expander, where $\beta_n \gg \ln n$. Define the following quantities:*

$$\begin{aligned} \epsilon &= \frac{1-\rho}{2}, \\ \hat{\rho} &= \frac{1}{1+\epsilon/4}, \\ b_n &= \frac{320}{(1-\rho)^2} \cdot \frac{n \ln n}{\beta_n} = \frac{80}{\epsilon^2} \cdot \frac{n \ln n}{\beta_n}, \\ m_n &= \frac{\rho}{\hat{\rho}} b_n, \\ \hat{u}_n &= \beta_n \frac{m_n}{n}. \end{aligned} \tag{21}$$

Let $\hat{\Gamma}$ be an arbitrary subset of the left vertices, I , such that

$$|\hat{\Gamma}| = m_n, \tag{22}$$

and let Δ be a random subset of the right vertices, J , where each vertex belongs to Δ independently and with the same probability, where

$$\mathbb{P}(j \in \Delta) \geq (\rho + 3\epsilon/4) \frac{b_n}{n}, \quad \forall j \in J, \quad (23)$$

for all n sufficiently large. Denote by \hat{G} the random subgraph $g_n|_{\hat{\Gamma} \cup \Delta}$. Then

$$\mathbb{P}\left(\hat{G} \text{ is not a } (\gamma/\hat{u}_n, \hat{u}_n)\text{-expander}\right) \leq \frac{1}{n^3}, \quad (24)$$

for all n sufficiently large, where the probability is measured with respect to the randomness in Δ .

To invoke Lemma 4.7, note that the conditions in Eq. (21) are identical to the definitions for the corresponding quantities in this section. We next verify that Eq. (23) is satisfied by the random subset, Δ , consisting of the idle servers at the end of a service slot. Recall that the length of a service slot is $\frac{b_n}{n}(\rho + \epsilon)$, and hence the probability that a given server, j , becomes idle by the end of a service slot is

$$\mathbb{P}(j \in \Delta) = 1 - \exp\left(-\frac{b_n}{n}(\rho + \epsilon)\right) \sim (\rho + \epsilon) \frac{b_n}{n}, \quad (25)$$

as $n \rightarrow \infty$. Therefore, for all n sufficiently large, we have that $\mathbb{P}(j \in \Delta) \geq (\rho + 3\epsilon/4) \frac{b_n}{n}$. We will now apply Lemmas 4.6 and 4.7 to the random subgraph with left (respectively, right) nodes $\hat{\Gamma}$ (respectively Δ), and with the demands A_i , for $i \in \hat{\Gamma}$, playing the role of λ .

Lemma 4.8 *If n is large enough, and if the value \mathbf{a} of \mathbf{A} is nice, then*

$$\mathbb{P}(X = 1 \mid \mathbf{A} = \mathbf{a}) \leq \frac{1}{n^3},$$

where the probability is with respect to the randomness in Δ .

Proof. We fix some \mathbf{a} , assumed to be nice. Recall that

$$\hat{u}_n = \beta_n \frac{m_n}{n},$$

and from the statement of Theorem 3.4 that

$$\gamma = \sqrt{\hat{\rho}} > \hat{\rho}.$$

We apply Lemma 4.6 to the randomly sampled subgraph \hat{G} , with left nodes $\hat{\Gamma}$, $|\hat{\Gamma}| = m_n$, and right nodes Δ . We have the following correspondence: the parameters n and ρ , in Lemma 4.6 become, in the current context, m_n and $\hat{\rho}$, respectively, and the parameters β_n and u_n both become \hat{u}_n . Thus, by Lemma 4.6,

$$\text{if } \hat{G} \text{ is a } (\gamma/\hat{u}_n, \hat{u}_n)\text{-expander, then } \mathbf{A}_{m_n}(\hat{u}_n) \subset \overline{\mathbf{R}}(\hat{G}). \quad (26)$$

Let $\hat{\mathbf{A}}$ be the vector of job arrival numbers \mathbf{A} , restricted to the set of nodes in $\hat{\Gamma}$, and let $\hat{\mathbf{a}}$ be the realization of $\hat{\mathbf{A}}$. Note that we have

$$\sum_{i \in \hat{\Gamma}} \hat{a}_i = \sum_{i=1}^n A_i = \rho b_n = \hat{\rho} m_n,$$

because of Eq. (20). Furthermore, for any $i \in \hat{\Gamma}$, the fact that \mathbf{a} is nice implies that $\hat{a}_i < \hat{u}_n$. Thus, $\hat{\mathbf{a}} \in \mathbf{L}_{m_n}(\hat{u}_n)$. By Eq. (26), this further implies that

$$\text{if } \hat{G} \text{ is a } (\gamma/\hat{u}_n, \hat{u}_n)\text{-expander, then } \hat{\mathbf{a}} \in \overline{\mathbf{R}}(\hat{G}). \quad (27)$$

By Lemma 4.7, the graph \hat{G} is a $(\gamma/\hat{u}_n, \hat{u}_n)$ -expander with probability at least $1 - n^{-3}$. Combining this fact with Eq. (27), we have thus verified that $\hat{\mathbf{a}}$ belongs to $\overline{\mathbf{R}}(\hat{G})$, with probability at least

$$1 - \frac{1}{n^3}.$$

With $\overline{\mathbf{R}}(\hat{G})$ having been defined as the closure of the capacity region, $\mathbf{R}(\hat{G})$ (cf. Definition 2.2), the fact that the vector $\hat{\mathbf{a}}$ belongs to $\overline{\mathbf{R}}(\hat{G})$ is a statement about the existence of a feasible flow, $\{\hat{f}_{ij} : (i, j) \in \hat{E}\}$ (where \hat{E} is the set of edges in \hat{G}), in a linear network flow model of the form

$$\begin{aligned} \hat{a}_i &= \sum_{j: (i,j) \in \hat{E}} f_{ij}, \quad \forall i \in \hat{\Gamma}, \\ \sum_{i: (i,j) \in \hat{E}} f_{ij} &\leq 1, \quad \forall j \in \Delta, \\ f_{ij} &\geq 0, \quad \forall (i, j) \in \hat{E}. \end{aligned}$$

Because the ‘‘supplies’’ \hat{a}_i in this network flow model, as well as the unit capacities of the right nodes are integer, it is well known that there also exists an integer flow. That is, we can find $\hat{f}_{ij} \in \{0, 1\}$ such that $\sum_j \hat{f}_{ij} = \hat{a}_i$, for all i , and $\sum_i \hat{f}_{ij} \leq 1$, for all j . But this is the same as the statement that there exists a feasible batch assignment over \hat{G} . Thus, for large enough n and for any given nice \mathbf{a} , the conditional probability that a batch assignment does not exist is upper bounded by n^{-3} , as claimed. Q.E.D.

We can now complete the proof of Proposition 4.4. By considering unconditional probabilities where \mathbf{A} is random, and for n large enough, we have that

$$\begin{aligned} \mathbb{P}(X = 1) &\leq \mathbb{P}(\mathbf{A} \text{ is not nice}) + \sum_{\mathbf{a} \text{ nice}} \mathbb{P}(X = 1 \mid \mathbf{A} = \mathbf{a}) \cdot \mathbb{P}(\mathbf{A} = \mathbf{a}) \\ &\stackrel{(a)}{\leq} \frac{1}{n^3} + \sum_{\mathbf{a} \text{ nice}} \mathbb{P}(X = 1 \mid \mathbf{A} = \mathbf{a}) \cdot \mathbb{P}(\mathbf{A} = \mathbf{a}) \\ &\stackrel{(b)}{\leq} \frac{1}{n^3} + \sum_{\mathbf{a} \text{ nice}} \frac{1}{n^3} \cdot \mathbb{P}(\mathbf{A} = \mathbf{a}) \\ &\leq \frac{2}{n^3} \\ &\leq \frac{1}{n^2}, \end{aligned} \quad (28)$$

where steps (a) and (b) follow from Lemmas 4.5 and 4.8, respectively. This concludes the proof of Proposition 4.4.

4.8. Service and Waiting Time Bounds for the Virtual Queue

4.8.1. Service Time Bounds. We will now use Lemma 4.3 and Proposition 4.4 to bound the mean and variance of the service times in the modified virtual queue.

Lemma 4.9 *The modified batch service times, S'_k , are i.i.d., with*

$$\mathbb{E}(S'_k | g_n) \sim (\rho + \epsilon) \cdot \frac{b_n}{n}, \quad \text{and} \quad \text{Var}(S'_k | g_n) \lesssim \frac{b_n^2}{n^2}.$$

Proof. We use the fact from Lemma 4.3, that $s \leq \mathbb{E}(S'_k | g_n) \leq s + q(g_n)c_n$, where $c_n \lesssim b_n$. We recall that $s = (\rho + \epsilon)b_n/n$, and use the fact $q(g_n) \leq n^{-2}$, as guaranteed by Proposition 4.4. The term $q(g_n)c_n$ satisfies $q(g_n)c_n \lesssim b_n/n^2$, which is of lower order than b_n/n , and hence negligible compared to s . This proves the first part of the lemma.

For the second part, we use Lemma 4.3 in the first inequality below, and the fact that $q(g_n) \leq n^{-2}$ in the second, to obtain

$$\text{Var}(S'_k | g_n) \lesssim q(g_n)c_n^2 \lesssim \frac{b_n^2}{n^2}.$$

Q.E.D.

4.8.2. Waiting Time Bounds. Fix n and the graph g_n . Let W^B be a random variable whose distribution is the same as the steady-state distribution of the time that a batch spends waiting in the queue of the virtual queueing system introduced in Section 4.3.2.

Proposition 4.10 *We have that*

$$\mathbb{E}(W^B | g_n) \lesssim \frac{b_n}{n}. \tag{29}$$

Proof. As discussed in Section 4.4, the waiting time of a batch, in the virtual queueing system, is dominated by the waiting time in a modified virtual queueing system, which is a GI/GI/1 queue, with independent inter-arrival times A_k (defined in Section 4.3.1) and independent service times S'_k . Let W' be a random variable whose distribution is the same as the steady-state distribution of the time that a batch spends waiting in the queue of the modified virtual queueing system.

According to Kingman's bound [18], W' satisfies

$$\mathbb{E}(W' | g_n) \leq \tilde{\lambda} \frac{\sigma_a^2 + \sigma_s^2}{2(1 - \tilde{\rho})},$$

where $\tilde{\lambda}$ is the arrival rate, $\tilde{\rho}$ is the traffic intensity, and σ_a^2 and σ_s^2 are the variances of the inter-arrival times and service times, respectively, that are associated with the modified virtual queueing system.

From Lemma 4.2, we have

$$\tilde{\lambda} = \frac{1}{\mathbb{E}(A_k)} \leq \frac{n}{b_n}.$$

and

$$\sigma_a^2 = \text{Var}(A_k) \lesssim \frac{b_n}{n^2}.$$

We now bound

$$\tilde{\rho} = \frac{\mathbb{E}(S'_k | g_n)}{\mathbb{E}(A_k)}.$$

From the first part of Lemma 4.9, we have $\mathbb{E}(S'_k | g_n) \sim (\rho + \epsilon)b_n/n$. Together with the bound $1/\mathbb{E}(A_k) \leq n/b_n$, we obtain that as $n \rightarrow \infty$, $\tilde{\rho}$ is upper bounded by a number strictly less than 1.

We also have, from the second part of Lemma 4.9,

$$\sigma_s^2 = \text{Var}(S'_k | g_n) \lesssim \frac{b_n^2}{n^2}.$$

Using these inequalities in Kingman's bound, we obtain

$$\mathbb{E}(W^B | g_n) \leq \mathbb{E}(W' | g_n) \lesssim \frac{n}{b_n} \cdot \frac{b_n^2}{n^2} = \frac{b_n}{n}.$$

Q.E.D.

4.9. Completing the Proof of Theorem 3.4

Proof. As discussed in Section 4.3.2, the expected waiting time of a job is upper bounded by the sum of three quantities.

(a) The expected time from the arrival of the job until the arrival time of the batch that the job belongs to. This is bounded above by the expected time until there are ρb_n subsequent arrivals, which is equal to $\mathbb{E}(A_1)$. By Lemma 4.2, this is bounded above by $c_1 b_n/n$, for some constant c_1 .

(b) The expected time that the batch waits in the virtual queue. This is also upper bounded by $c_2 b_n/n$, by Proposition 4.10, for some constant c_2 .

(c) The service time of the batch, which (by Lemma 4.9) again admits an upper bound of the form $c_3 b_n/n$, for some constant c_3 .

Furthermore, in the results that give these upper bounds, c_1 , c_2 , and c_3 , are absolute constants, that do not depend on λ_n or g_n .

By our assumptions on the choice of b_n in Section 4.6, we have $b_n = \frac{320}{(1-\rho)^2} \cdot \frac{n \ln n}{\beta_n}$, and β_n is proportional to d_n . We conclude that there exists a constant c such that for large enough n , we have $\mathbb{E}(W | g_n, \lambda_n) \leq c \ln n/d_n$, for any given $\lambda_n \in \Lambda_n(u_n)$, which is an upper bound of the desired form. This establishes Part 1 of the theorem. Finally, Part 2 follows from the way that the policy was constructed. Q.E.D.

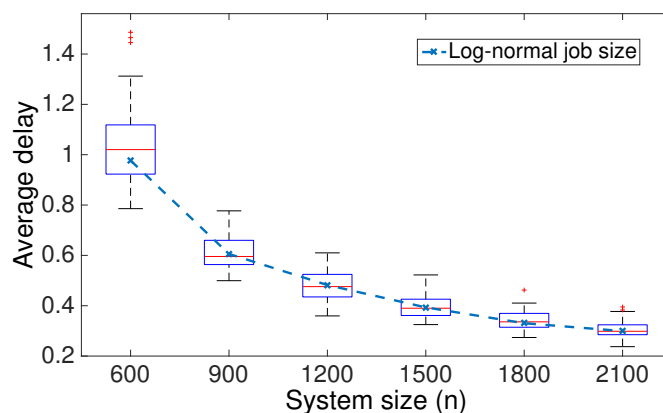


Figure 5 Simulations of the virtual-queue based policy given in Section 4.3, with $d_n = n^{2/3}$, $b_n = n \ln(n)/d_n$, and $\lambda_i = 0.5$ for all $i = 1, \dots, n$. The boxplot contains the average delay from 50 runs of simulations where the job size distribution is assumed to be exponential with mean 1. Each run is performed on a random d_n -regular graph over 10^4 service slots and a 1000-slot burn-in period. The center line of a box represents the median and upper and lower edges of the box represent the 25th and 75th percentiles, respectively. The dashed line depicts the median average waiting times when the job sizes are distributed according to a log-normal distribution with mean 1 and variance 10.

4.10. On Practical Policies

Figure 5 provides simulation results for the average delay under the virtual-queue based scheduling policy used in proving Theorem 3.4. The main role of the policy is to demonstrate the fundamental potential of the Expander architecture in jointly achieving a small delay and large capacity region when the system size is large. In smaller systems, however, there could be other policies that yield better performance. For instance, simulations suggest that a seemingly naive greedy heuristic can achieve a smaller delay in moderately-sized systems, which is practically zero in the range of parameters in Figure 5. Under the greedy heuristic, an available server simply fetches a job from a longest connected queue, and a job is immediately sent to a connected idle server upon arrival if possible. Intuitively, the greedy policy can provide a better delay because it avoids the overhead of holding jobs in queues while forming a batch. Unfortunately, it appears challenging to establish rigorous delay or capacity guarantees for the greedy heuristic and other similar policies.

In some applications, such as call centers, the service times or job sizes may not be exponentially distributed ([4]). In Figure 5, we also include the scenario where the job sizes are drawn from a log-normal distribution ([4]) with an increased variance. Interestingly, the average delay appears to be somewhat insensitive to the change in job size distribution.

5. Summary and Future Research

The main message of this paper is that the two objectives of a large capacity region and an asymptotically vanishing delay can be simultaneously achieved even if the level of processing flexibility of each server is small compared to the system size. Our main results show that, as far as these objectives are concerned, the family of Expander architectures is essentially optimal: it admits a capacity region whose size is within a constant factor of the maximum possible, while ensuring an asymptotically vanishing queueing delay for all arrival rate vectors in the capacity region.

An alternative design, the Random Modular architecture, guarantees small delays for “many” arrival rates, by means of a simple greedy scheduling policy. However, for any given Modular architecture, there are always many arrival rate vectors in $\mathbf{\Lambda}_n(u_n)$ that result in an unstable system, even if the maximum arrival rate across the queues is of constant order. Nevertheless, the simplicity of the Modular architectures can still be appealing in some practical settings.

Our result for the Expander architecture leaves open three questions:

1. Is it possible to lower the requirement on the average degree from $d_n \gg \ln n$ to $d_n \gg 1$?
2. Without sacrificing the size of the capacity region, is it possible to achieve a queueing delay which approaches zero exponentially fast as a function of d_n ? The delay scaling in Theorem 3.4 is $\mathcal{O}(\ln n/d_n)$.
3. Is it possible to obtain delay and stability guarantees under simpler policies, such as the greedy heuristic mentioned in Section 4.10? The techniques developed in [31] for analyzing first-come-first-serve scheduling rules in a multi-class queueing network similar to ours could be a useful starting point.

Finally, the scaling regime considered in this paper assumes that the traffic intensity is fixed as n increases, which fails to capture system performance in the heavy-traffic regime ($\rho \approx 1$). It would be interesting to consider a scaling regime in which ρ and n scale simultaneously (e.g., as in the celebrated Halfin-Whitt regime [11]), but it is unclear at this stage what the most appropriate formulations and analytical techniques are.

References

- [1] A. S. Asratian, T. M. J. Denley, and R. Haggkvist. *Bipartite Graphs and their Applications*. Cambridge University Press, 1998.
- [2] A. Bassamboo, R. S. Randhawa, and J. A. V. Mieghem. A little flexibility is all you need: on the asymptotic value of flexible capacity in parallel queueing systems. *Operations Research*, 60(6):1423–1435, December 2012.
- [3] S. L. Bell and R. J. Williams. Dynamic scheduling of a system with two parallel servers in heavy traffic with resource pooling: asymptotic optimality of a threshold policy. *Ann. Appl. Probab.*, 11(3):608–649, 2001.

-
- [4] L. Brown, N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn, and L. Zhao. Statistical analysis of a telephone call center: A queueing-science perspective. *Journal of the American statistical association*, 100(469):36–50, 2005.
- [5] X. Chen, J. Zhang, and Y. Zhou. Optimal sparse designs for process flexibility via probabilistic expanders. *Operations Research*, 63(5):1159–1176, 2015.
- [6] M. Chou, G. A. Chua, C.-P. Teo, and H. Zheng. Design for process flexibility: efficiency of the long chain and sparse structure. *Operations Research*, 58(1):43–58, 2010.
- [7] M. Chou, C.-P. Teo, and H. Zheng. Process flexibility revisited: the graph expander and its applications. *Operations Research*, 59(5):1090–1105, 2011.
- [8] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.
- [9] D. Gross, J. F. Shortle, J. M. Thompson, and C. M. Harris. *Fundamentals of Queueing Theory*. John Wiley & Sons, 2008.
- [10] S. Gurumurthi and S. Benjaafar. Modeling and analysis of flexible queueing systems. *Management Science*, 49:289–328, 2003.
- [11] S. Halfin and W. Whitt. Heavy-traffic limits for queues with many exponential servers. *Operations Research*, 29:567–588, 1981.
- [12] J. M. Harrison and M. J. Lopez. Heavy traffic resource pooling in parallel-server systems. *Queueing Systems*, 33:39–368, 1999.
- [13] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- [14] S. Hoory, N. Linial, and A. Wigderson. Expander graphs and their applications. *Bulletin of the American Mathematical Society*, 43(4):439–561, 2006.
- [15] S. M. Iravani, M. P. V. Oyen, and K. T. Sims. Structural flexibility: A new perspective on the design of manufacturing and service operations. *Management Science*, 51(2):151–166, 2005.
- [16] W. Jordan and S. C. Graves. Principles on the benefits of manufacturing process flexibility. *Management Science*, 41(4):577–594, 1995.
- [17] S. Kandula, S. Sengupta, A. Greenberg, P. Patel, and R. Chaiken. The nature of data center traffic: measurements & analysis. In *ACM SIGCOMM*, 2009.
- [18] J. Kingman. Some inequalities for the queue GI/G/1. *Biometrika*, 49(3/4):315–324, 1962.
- [19] S. Kunniyur and R. Srikant. Analysis and design of an adaptive virtual queue. In *ACM SIGCOMM*, 2001.
- [20] M. Leconte, M. Lelarge, and L. Massoulié. Bipartite graph structures for efficient balancing of heterogeneous loads. *ACM SIGMETRICS Performance Evaluation Review*, 40(1):41–52, 2012.

- [21] A. Mandelbaum and M. I. Reiman. On pooling in queueing networks. *Management Science*, 44(7):971–981, 1998.
- [22] N. McKeown, A. Mekkittikul, V. Anantharam, and J. Walrand. Achieving 100% throughput in an input-queued switch. *IEEE Trans. on Comm*, 47(8):1260–1267, 1999.
- [23] M. Mitzenmacher and E. Upfal. *Probability and computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, 2005.
- [24] M. Neely, E. Modiano, and Y. Cheng. Logarithmic delay for $n \times n$ packet switches under the crossbar constraint. *IEEE/ACM Trans. Netw*, 15(3):657–668, 2007.
- [25] D. Shah and J. N. Tsitsiklis. Bin packing with queues. *Journal of Applied Probability*, 45(4):922–939, 2008.
- [26] D. Simchi-Levi and Y. Wei. Understanding the performance of the long chain and sparse designs in process flexibility. *Operations Research*, 60(5):1125–1141, 2012.
- [27] G. Soundararajan, C. Amza, and A. Goel. Database replication policies for dynamic content applications. In *ACM SIGOPS Operating Systems Review*, volume 40, pages 89–102. ACM, 2006.
- [28] R. Talreja and W. Whitt. Fluid models for overloaded multiclass many-server queueing systems with first-come, first-served routing. *Management Science*, 54(8):1513–1527, 2008.
- [29] J. N. Tsitsiklis and K. Xu. On the power of (even a little) resource pooling. *Stochastic Systems*, 2(1):1–66, 2012. Available at <http://dx.doi.org/10.1214/11-SSY033>.
- [30] J. N. Tsitsiklis and K. Xu. Queueing system topologies with limited flexibility. *ACM SIGMETRICS Performance Evaluation Review*, 41(1):167–178, 2013.
- [31] J. Visschers, I. Adan, and G. Weiss. A product form solution to a system with multi-type jobs and multi-type servers. *Queueing Systems*, 70:269–298, 2012.
- [32] R. Wallace and W. Whitt. A staffing algorithm for call centers with skill-based routing. *Manufacturing and Service Operations Management*, 7:276–294, 2005.
- [33] K. Xu. *On the power of (even a little) flexibility in dynamic resource allocation*. PhD thesis, Massachusetts Institute of Technology, 2014. Available at <http://hdl.handle.net/1721.1/91101>.

Appendix A: Proofs

A.1. Proof of Lemma 3.2

Proof. Fix $\lambda = (\lambda_1, \dots, \lambda_n) \in \Lambda_n(u_n)$, and let g_n be a $(\gamma/\beta_n, \beta_n)$ -expander, where $\gamma > \rho$ and $\beta_n \geq u_n$. By the max-flow-min-cut theorem, and the fact that all servers have unit capacity, it suffices to show that

$$\sum_{i \in S} \lambda_i < |\mathcal{N}(S)|, \quad \forall S \subset I. \quad (30)$$

We consider two cases, depending on the size of S .

1. Suppose that $|S| < \gamma n / \beta_n$. By the expansion property of g_n , we have that

$$\mathcal{N}(S) \geq \beta_n |S| \geq u_n |S| > \sum_{i \in S} \lambda_i, \quad (31)$$

where the second inequality follows from the fact that $\beta_n \geq u_n$, and the last inequality from $\lambda_i < u_n$ for all $i \in I$.

2. Suppose that $|S| \geq \gamma n / \beta_n$. By removing, if necessary, some of the nodes in S , we obtain a set $S' \subset S$ of size exactly $\gamma n / \beta_n$, and

$$\mathcal{N}(S) \geq \mathcal{N}(S') \stackrel{(a)}{\geq} \gamma n > \rho n \stackrel{(b)}{\geq} \sum_{i \in S} \lambda_i, \quad (32)$$

where step (a) follows from the expansion property, and step (b) from the assumption that $\sum_{i \in I} \lambda_i \leq \rho n$.

This completes the proof. Q.E.D.

A.2. Proof of Lemma 3.3

Proof. Lemma 3.3 is a consequence of the following standard result (cf. [1]), where we let $d = d_n$, $\beta = \beta_n$, and $\alpha = \gamma / \beta_n = \sqrt{\rho} / \beta_n$, and observe that $\log_2 \beta_n \ll \beta_n$ as $n \rightarrow \infty$.

Lemma A.1 *Fix $n \geq 1$, $\beta \geq 1$ and $\alpha\beta < 1$. If*

$$d \geq \frac{1 + \log_2 \beta + (\beta + 1) \log_2 e}{-\log_2(\alpha\beta)} + \beta + 1, \quad (33)$$

then there exists an (α, β) -expander with maximum degree d .

Q.E.D.

A.3. Proof of Theorem 3.5

Proof. Since the arrival rate vector λ_n whose existence we want to show can depend on the architecture, we assume, without loss of generality, that servers and queues are clustered in the same manner: server i and queue i belong to the same cluster. Since all servers have capacity 1, and each cluster has exactly d_n servers, it suffices to show that there exists $\lambda = (\lambda_1, \dots, \lambda_n) \in \mathbf{\Lambda}_n(u_n)$, such that the total arrival rate to the first queue cluster exceeds d_n , i.e.,

$$\sum_{i=1}^{d_n} \lambda_i > d_n. \quad (34)$$

To this end, consider the vector λ where $\lambda_i = \min\{2, (1 + u_n)/2\}$ for all $i \in \{1, \dots, d_n\}$, and $\lambda_i = 0$ for $i \geq d_n + 1$. Because of the assumption $u_n > 1$ in the statement of the theorem, we have that

$$\max_{1 \leq i \leq n} \lambda_i = \min\{2, (1 + u_n)/2\} \leq \frac{1 + u_n}{2} < u_n, \quad (35)$$

and

$$\sum_{i=1}^n \lambda_i = d_n \min\{2, (1 + u_n)/2\} \leq 2d_n \leq 2 \cdot \frac{\rho}{2} n = \rho n, \quad (36)$$

where the last inequality in Eq. (36) follows from the assumption that $d_n \leq \frac{\rho}{2} n$. Eqs. (35) and (36) together ensure that $\boldsymbol{\lambda} \in \boldsymbol{\Lambda}_n(u_n)$ (cf. Condition 1). Since we have assumed that $u_n > 1$, we have $\lambda_i > 1$, for $i = 1, \dots, d_n$, and therefore Eq. (34) holds for this $\boldsymbol{\lambda}$. We thus have that $\boldsymbol{\lambda} \notin \mathbf{R}(g_n)$, which proves our claim. Q.E.D.

A.4. Proof of Theorem 3.6

Proof. Part (a); Eq. (5). We will use the following classical result due to Hoeffding, adapted from Theorem 3 in [13].

Lemma A.2 *Fix integers m and n , where $0 < m < n$. Let X_1, X_2, \dots, X_m be random variables drawn uniformly from a finite set $C = \{c_1, \dots, c_n\}$, without replacement. Suppose that $0 \leq c_i \leq b$ for all i , and let $\sigma^2 = \text{Var}(X_1)$. Let $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$. Then,*

$$\mathbb{P}(\bar{X} \geq \mathbb{E}(\bar{X}) + t) \leq \exp\left(-\frac{mt}{b} \left[\left(1 + \frac{\sigma^2}{bt}\right) \ln\left(1 + \frac{bt}{\sigma^2}\right) - 1\right]\right), \quad (37)$$

for all $t \in (0, b)$.

We fix some $\boldsymbol{\lambda}_n \in \boldsymbol{\Lambda}_n(u_n)$. If $u_n < 1$, then $\boldsymbol{\lambda}_n \in \boldsymbol{\Lambda}_n(1)$. It therefore suffices to prove the result for the case where $u_n \geq 1$ and we will henceforth assume that this is the case. Recall that $A_k \subset I$ is the set of d_n queues in the k th queue cluster generated by the partition $\sigma_n = (A_1, \dots, A_{n/d_n})$. We consider some $\epsilon \in (0, 1/\rho)$, and define the event E_k as

$$E_k = \left\{ \sum_{i \in A_k} \lambda_i > (1 + \epsilon)\rho d_n \right\}. \quad (38)$$

Since σ_n is drawn uniformly at random from all possible partitions, it is not difficult to see that $\sum_{i \in A_k} \lambda_i$ has the same distribution as $\sum_{i=1}^{d_n} X_i$, where X_1, X_2, \dots, X_{d_n} are d_n random variables drawn uniformly at random, without replacement, from the set $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$. Note that $\epsilon\rho < 1 \leq u_n$, so that $\epsilon\rho \in (0, u_n)$. We can therefore apply Lemma A.2, with $m = d_n$, $b = u_n$, and $t = \epsilon\rho$, to obtain

$$\begin{aligned} \mathbb{P}(E_1) &= \mathbb{P}\left(\sum_{i=1}^{d_n} X_i > (1 + \epsilon)\rho d_n\right) \\ &\stackrel{(a)}{\leq} \mathbb{P}\left(\frac{1}{d_n} \sum_{i=1}^{d_n} X_i > \mathbb{E}\left(\frac{1}{d_n} \sum_{i=1}^{d_n} X_i\right) + \epsilon\rho\right) \\ &\leq \exp\left(-\frac{\epsilon\rho d_n}{u_n} \left[\left(1 + \frac{\text{Var}(X_1)}{\epsilon\rho u_n}\right) \ln\left(1 + \frac{\epsilon\rho u_n}{\text{Var}(X_1)}\right) - 1\right]\right), \end{aligned} \quad (39)$$

where the probability is taken with respect to the randomness in G , and where in step (a) we used the fact that

$$\mathbb{E} \left(\sum_{i=1}^{d_n} X_i \right) = \sum_{i=1}^{d_n} \mathbb{E}(X_i) = d_n \mathbb{E}(X_1) = d_n \left(\frac{1}{n} \sum_{i=1}^n \lambda_i \right) \leq \rho d_n. \quad (40)$$

We now develop an upper bound on $\text{Var}(X_1)$. Since X_1 takes values in $[0, u_n]$, we have $X_1^2 \leq u_n X_1$ and, therefore,

$$\text{Var}(X_1) \leq \mathbb{E}(X_1^2) \leq u_n \mathbb{E}(X_1) \leq \rho u_n. \quad (41)$$

Observe that for all $a, x > 0$,

$$\frac{d}{dx} (1 + x/a) \ln(1 + a/x) = -\frac{1}{x} + \frac{1}{a} \ln(1 + a/x) < -\frac{1}{x} + \frac{1}{a} \cdot \frac{a}{x} = 0. \quad (42)$$

Therefore, with the substitutions $a = \epsilon \rho u_n$ and $x = \text{Var}(X_1)$, we have that the right-hand-side of (39) is increasing in $\text{Var}(X_1)$. Combining Eqs. (39) and (41), we obtain

$$\mathbb{P}(E_1) \leq \exp \left(-\frac{\epsilon \rho d_n}{u_n} \left[\left(1 + \frac{1}{\epsilon} \right) \ln(1 + \epsilon) - 1 \right] \right).$$

Note that

$$\frac{d}{dx} \left(1 + \frac{1}{x} \right) \ln(1 + x) = \frac{1}{x^2} (x - \ln(1 + x)) \xrightarrow{(a)} \frac{1}{2}, \quad \text{as } x \downarrow 0, \quad (43)$$

where step (a) follows from applying l'Hôpital's rule. We thus have that $[(1 + \frac{1}{\epsilon}) \ln(1 + \epsilon) - 1] \sim \frac{1}{2} \epsilon \geq \frac{1}{3} \epsilon$, as $\epsilon \downarrow 0$, it follows that there exists $\theta > 0$ such that for all $\epsilon \in (0, \theta)$,

$$\mathbb{P}(E_1) \leq \exp \left(-\frac{\rho}{3} \cdot \frac{\epsilon^2 d_n}{u_n} \right). \quad (44)$$

Let $\epsilon = \frac{1}{2} \min\{\frac{1}{\rho} - 1, \theta\}$; in particular, our earlier assumption that $\epsilon \rho < 1$ is satisfied. Suppose that $u_n \leq \frac{\rho \epsilon^2}{6} d_n \ln^{-1} n$. Combining Eq. (44) with the union bound, we have that

$$\begin{aligned} \mathbb{P}_{G_n}(\boldsymbol{\lambda}_n \notin \mathbf{R}(G_n)) &\leq \mathbb{P} \left(\bigcup_{k=1}^{n/d_n} E_k \right) \\ &\leq \sum_{k=1}^{n/d_n} \mathbb{P}(E_k) \\ &\leq \frac{n}{d_n} \exp \left(-\frac{\rho}{3} \cdot \frac{\epsilon^2 d_n}{u_n} \right) \\ &\stackrel{(a)}{\leq} \frac{n}{d_n} \cdot \frac{1}{n^2} \\ &\leq n^{-1}, \end{aligned} \quad (45)$$

where step (a) follows from the assumption that $u_n \leq \frac{\rho \epsilon^2}{6} d_n \ln^{-1} n$. It follows that

$$\lim_{n \rightarrow \infty} \inf_{\boldsymbol{\lambda}_n \in \boldsymbol{\Lambda}_n(u_n)} \mathbb{P}_{G_n}(\boldsymbol{\lambda}_n \in \mathbf{R}(G_n)) \geq \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n} \right) = 1. \quad (46)$$

We have therefore proved part (a) of the theorem, with $c_2 = \rho\epsilon^2/6$.

Part (b); Eq. (6).

Let us fix a large enough constant c_3 , whose value will be specified later, and let

$$v_n = c_3 \frac{d_n}{\ln n}. \quad (47)$$

For this part of the proof, we will assume that $u_n > v_n$. Because we are interested in showing a result for the worst case over all $\lambda_n \in \Lambda_n(u_n)$, we can assume that $u_n \ll n$.

At this point, we could analyze the model for a worst-case choice of λ_n . However, the analysis turns out to be simpler if we employ the probabilistic method. Denote by μ_n a probability measure over $\Lambda_n(u_n)$. Let λ_n be a random vector drawn from the distribution μ_n , independent of the randomness in the Random Modular architecture, G . (For convenience, we suppress the subscript n and write G instead of G_n .) The following elementary fact captures the essence of the probabilistic method.

Lemma A.3 *Fix n , a measure μ_n on $\Lambda_n(u_n)$, and a constant a_n . Suppose that*

$$\mathbb{P}_{\lambda_n, G}(\lambda_n \notin \mathbf{R}(G)) \geq a_n, \quad (48)$$

where $\mathbb{P}_{\lambda_n, G}$ stands for the product of the measures μ_n (for λ_n) and \mathbb{P}_G (for G). Then,

$$\sup_{\tilde{\lambda}_n \in \Lambda_n(u_n)} \mathbb{P}_G(\tilde{\lambda}_n \notin \mathbf{R}(G)) \geq a_n. \quad (49)$$

Proof. We have that

$$\begin{aligned} \sup_{\tilde{\lambda}_n \in \Lambda_n(u_n)} \mathbb{P}_G(\tilde{\lambda}_n \notin \mathbf{R}(G)) &\geq \int_{\tilde{\lambda}_n \in \Lambda_n(u_n)} \mathbb{P}_G(\tilde{\lambda}_n \notin \mathbf{R}(G)) d\mu_n(\tilde{\lambda}_n) \\ &= \mathbb{P}_{\lambda_n, G}(\lambda_n \notin \mathbf{R}(G)) \\ &\geq a_n. \end{aligned} \quad (50)$$

Q.E.D.

We will now construct sequences, $\{\mu_n : n \in \mathbb{N}\}$, and $\{a_n : n \in \mathbb{N}\}$, with $\lim_{n \rightarrow \infty} a_n = 1$, so that Eq. (48) holds for all n . To simplify notation, in the rest of this proof we will write \mathbb{P} instead of \mathbb{P}_G or $\mathbb{P}_{\lambda_n, G}$, etc. Which particular measure we are dealing with will always be clear from the context.

Fix $n \in \mathbb{N}$. We first construct the distribution μ_n . Let $\lambda' = (\lambda'_1, \lambda'_2, \dots, \lambda'_n)$ be a random vector with independent components and with

$$\lambda'_i = \begin{cases} v_n, & \text{w.p. } \frac{\rho}{(1+\epsilon)v_n}, \\ 0, & \text{otherwise,} \end{cases} \quad (51)$$

for all i . Let H be the event defined by

$$H = \left\{ \sum_{i=1}^n \lambda'_i \leq \rho n \right\}. \quad (52)$$

Let λ_n be the random vector given by

$$\lambda_n = \mathbb{I}(H)\lambda', \quad (53)$$

where $\mathbf{0}$ is the zero vector of dimension n , and where $\mathbb{I}(\cdot)$ is the indicator function. That is, λ_n takes on the value of λ' if H occurs, and is set to zero, otherwise. It is not difficult to verify that, by construction, we always have $\lambda_n \in \Lambda_n(u_n)$. We let μ_n be the distribution of this random vector λ_n .

We next show that

$$\lim_{n \rightarrow \infty} \mathbb{P}(\lambda_n \notin \mathbf{R}(G)) = 1, \quad (54)$$

which, together with Lemma A.3 above, will complete the proof of the theorem. Fix some $\epsilon > \frac{1}{\rho} - 1$, so that $(1 + \epsilon)\rho > 1$, and define the event

$$E_k = \left\{ \sum_{i \in A_k} \lambda'_i > (1 + \epsilon)\rho d_n \right\}, \quad k \in \{1, \dots, n/d_n\}. \quad (55)$$

Note that, if some E_k occurs, then λ' will not be in $\mathbf{R}(G)$. Therefore,

$$\mathbb{P}(\lambda' \notin \mathbf{R}(G)) \geq \mathbb{P}\left(\bigcup_{k=1}^{n/d_n} E_k\right). \quad (56)$$

Let X_1, X_2, \dots be i.i.d. Bernoulli random variables with

$$\mathbb{E}(X_1) = \mathbb{P}(X_1 = 1) = \frac{\rho}{(1 + \epsilon)v_n}. \quad (57)$$

By the definition of λ' (cf. Eq. (51)), we have that

$$\begin{aligned} \mathbb{P}(E_1) &= \mathbb{P}\left(\sum_{i \in A_1} \lambda'_i > (1 + \epsilon)\rho d_n\right) \\ &= \mathbb{P}\left(\sum_{i=1}^{d_n} X_i > (1 + \epsilon)\rho \frac{d_n}{v_n}\right) \\ &= \mathbb{P}\left(\frac{1}{d_n} \sum_{i=1}^{d_n} X_i > (1 + \epsilon)^2 \mathbb{E}(X_1)\right). \end{aligned} \quad (58)$$

By Sanov's theorem (cf. Chapter 12 of [8]), we have that

$$\begin{aligned} \mathbb{P}(E_1) &= \mathbb{P}\left(\frac{1}{d_n} \sum_{i=1}^{d_n} X_i > (1 + \epsilon)^2 \mathbb{E}(X_1)\right) \\ &\gtrsim \frac{1}{d_n^2} \exp\left(-D_B\left(\frac{(1 + \epsilon)\rho}{v_n} \parallel \frac{\rho}{(1 + \epsilon)v_n}\right) d_n\right), \end{aligned} \quad (59)$$

where $D_B(p||q)$ is the Kullback-Leibler divergence between two Bernoulli distributions with parameters p and q , respectively:

$$D_B(p||q) = p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q}. \quad (60)$$

Let us fix some $r \in (0, 1)$. Using the fact that $\ln(1+y) \sim y$ as $y \rightarrow 0$, we have that

$$D_B(x||rx) \sim x \left[\ln \frac{1}{r} + (1-r) \right], \quad \text{as } x \rightarrow 0. \quad (61)$$

Recall that $d_n \geq c_1 \ln n$ and $v_n \geq / \ln n$. By Eq. (61), with $x = (1+\epsilon)\rho/v_n$, $r = 1/(1+\epsilon)^2$, and for the given c_1 , we can set c_3 to be sufficiently large so that

$$\begin{aligned} D_B \left(\frac{(1+\epsilon)\rho}{v_n} \parallel \frac{\rho}{(1+\epsilon)v_n} \right) &\leq 2 \frac{(1+\epsilon)\rho}{v_n} \cdot \left[\ln(1+\epsilon)^2 + \left(1 - \frac{1}{(1+\epsilon)^2} \right) \right] \\ &= \frac{2h}{v_n}, \end{aligned} \quad (62)$$

for all sufficiently large n , where $h = (1+\epsilon)\rho \left[\ln(1+\epsilon)^2 + \left(1 - \frac{1}{(1+\epsilon)^2} \right) \right] > 0$. Combining Eqs. (59) and (62), we have that

$$\mathbb{P}(E_1) \gtrsim \frac{1}{d_n^2} \exp \left(-2h \frac{d_n}{v_n} \right) \stackrel{(a)}{\gtrsim} \frac{1}{d_n^2} n^{-2h/c_3}, \quad (63)$$

where step (a) follows from the assumption that $v_n \geq c_3 d_n / \ln n$. Equation (63) can be rewritten in the form

$$\mathbb{P}(E_1) \geq \frac{c}{d_n^2} n^{-2h/c_3}, \quad (64)$$

where c is a positive constant, and where the inequality is valid for large enough n .

Fix $c_3 = 40h$, and recall that $\epsilon > \frac{1}{\rho} - 1$. We have that

$$\begin{aligned} \mathbb{P}(\boldsymbol{\lambda}' \notin \mathbf{R}(G)) &\geq \mathbb{P} \left(\bigcup_{k=1}^{n/d_n} E_k \right) \\ &\stackrel{(a)}{=} 1 - \prod_{k=1}^{n/d_n} (1 - \mathbb{P}(E_k)) \\ &= 1 - (1 - \mathbb{P}(E_1))^{n/d_n} \\ &\stackrel{(b)}{\geq} 1 - \left(1 - c d_n^{-3} n^{1-2h/c_3} d_n/n \right)^{n/d_n} \\ &\stackrel{(c)}{\geq} 1 - \left(1 - c n^{0.05} d_n/n \right)^{n/d_n} \\ &\rightarrow 1, \quad \text{as } n \rightarrow \infty, \end{aligned} \quad (65)$$

where step (a) is based on the independence among the events E_k , which is in turn based on the independence among the λ'_i s; step (b) follows from Eq. (64) and some rearrangement; step (c)

follows from the assumption in the statement of the theorem that $d_n \leq n^{0.3}$, and our choice of $c_3 = 40h$.

We next show that the event H occurs with high probability when n is large. Let, as before, the X_i s be i.i.d. Bernoulli random variables with $\mathbb{E}(X_1) = \frac{\rho}{v_n(1+\epsilon)}$. Then,

$$\begin{aligned} \mathbb{P}(H) &= \mathbb{P}\left(\sum_{i=1}^n \lambda'_i \leq \rho n\right) \\ &= \mathbb{P}\left(\sum_{i=1}^n X_i \leq \rho n/v_n\right) \\ &= \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n X_i \leq (1+\epsilon)\mathbb{E}(X_1)\right) \rightarrow 1, \quad \text{as } n \rightarrow \infty, \end{aligned} \quad (66)$$

by the weak law of large numbers.

We are now ready to prove Eq. (54). We have that

$$\begin{aligned} \mathbb{P}_{\lambda_n, G}(\lambda_n \notin \mathbf{R}(G)) &= \mathbb{P}_{\lambda', G}(\mathbb{I}(H)\lambda' \notin \mathbf{R}(G)) \\ &= \mathbb{P}_{\lambda', G}(H \cap \{\lambda' \notin \mathbf{R}(G)\}) \\ &\geq \mathbb{P}(H) + \mathbb{P}(\lambda' \notin \mathbf{R}(G)) - 1 \\ &\rightarrow 1, \quad \text{as } n \rightarrow \infty, \end{aligned} \quad (67)$$

where the last step follows from Eqs. (65) and (66). By Lemma A.3, Eq. (67) implies that $\lim_{n \rightarrow \infty} \sup_{\lambda_n \in \Lambda_n(u_n)} \mathbb{P}_{G_n}(\lambda_n \notin \mathbf{R}(G)) = 1$, which is in turn equivalent to $\lim_{n \rightarrow \infty} \inf_{\lambda_n \in \Lambda_n(u_n)} \mathbb{P}_{G_n}(\lambda_n \in \mathbf{R}(G)) = 0$. This proves Eq. (6). Q.E.D.

A.5. Proof of Theorem 3.7

Proof. Denote by $Q_i(t)$ the number of jobs in queue i at time t , and by $Q_k(t)$ the total number of jobs in queue cluster k , i.e.,

$$Q_k(t) = \sum_{i \in A_k} Q_i(t). \quad (68)$$

We note that $Q_k(\cdot)$ is the number of jobs in an $M/M/c$ queue, with $c = d_n$ and arrival rate $\eta_k = \sum_{i \in A_k} \lambda_i$. Also note that since $\lambda_n \in \gamma \mathbf{R}(g_n)$, we have that $\eta_k \leq \gamma d_n$. Using the formula for the expected waiting time in queue for an $M/M/c$ queue (cf. Section 2.3 of [9]), one can show that the average waiting time across jobs arriving to cluster k , W_k , satisfies

$$\mathbb{E}(W_k | \lambda) = \frac{1}{\sum_{i \in A_k} \lambda_i} \sum_{i \in A_k} \lambda_i \mathbb{E}(W_i) = \frac{C(d_n, \eta_k)}{d_n - \eta_k} \leq \frac{C(d_n, \gamma d_n)}{(1-\gamma)d_n} \lesssim \exp(-b \cdot d_n), \quad (69)$$

where $C(c, r)$ is given by

$$C(c, r) = \frac{r^c}{c!} \cdot \frac{1}{c(1-r/c)^2} \left(\frac{r^c}{c!} \cdot \frac{1}{1-r/c} + \sum_{i=0}^{c-1} \frac{r^i}{i!} \right)^{-1}.$$

The last inequality in Eq. (69) follows from the fact that for any given $\gamma \in (0, 1)$, there exists $b > 0$, so that $C(x, \gamma x) \lesssim \exp(-b \cdot x)$ as $x \rightarrow \infty$, as can be checked through elementary algebraic manipulations. Q.E.D.

A.6. Lower Bound on the Total Arrival Rate

We show in this section that the assumption that $\rho \in (1/2, 1)$ and $\sum_{i=1}^n \lambda_i \geq (1 - \rho)n$ (cf. Eq. (10) in Assumption 4.1) can be made without loss of generality. Fix the traffic intensity $\rho \in (0, 1)$, and suppose that $\boldsymbol{\lambda} \in \boldsymbol{\Lambda}_n(u_n)$. Define

$$\rho' = \rho + \frac{1}{2}(1 - \rho) = \frac{1 + \rho}{2}. \quad (70)$$

Note that $1/2 < \rho' < 1$, and $1 - \rho' = (1 - \rho)/2$. Consider a modified vector $\boldsymbol{\lambda}'$, where $\lambda'_i = (1 - \rho') + \lambda_i$, for all $i \in \{1, \dots, n\}$. By construction, we have that

$$\sum_{i=1}^n \lambda'_i \geq (1 - \rho')n, \quad (71)$$

$$\sum_{i=1}^n \lambda'_i \leq (1 - \rho')n + \sum_{i=1}^n \lambda_i \leq (1 - \rho')n + \rho n = \rho' n, \quad (72)$$

$$\max_{1 \leq i \leq n} \lambda'_i \leq \max_{1 \leq i \leq n} \lambda_i + (1 - \rho') < u_n + (1 - \rho'). \quad (73)$$

The above definition of $\boldsymbol{\lambda}'$ amounts to the following: we feed each queue with an additional independent Poisson stream of artificial (dummy) jobs of rate $1 - \rho'$. By Eqs. (72) and (73), the resulting arrival rate vector, $\boldsymbol{\lambda}'$, will belong to the set $\boldsymbol{\Lambda}_n(u_n + 1 - \rho')$. Also, by Eq. (71), it will satisfy the lower bound (10) on the total arrival rate, albeit with a modified traffic intensity of $\rho' \in (1/2, 1)$. Therefore, our assumption can always be satisfied by the insertion of dummy jobs. Note that the increment of $1 - \rho'$ to the value of u_n is insignificant in our regime of interest, where $u_n \gg 1$, and the insertion of dummy jobs only requires knowledge of the original traffic intensity, ρ .

A.7. Proof of Lemma 4.5

Proof. Note that because there are ρb_n jobs in a batch, the size of Γ is at most ρb_n , which is in turn less than m_n . This guarantees that the cardinality of $\hat{\Gamma}$ can be taken to be m_n . It therefore suffices to show that

$$\mathbb{P} \left(\max_{1 \leq i \leq n} A_i \geq \hat{u}_n \right) \leq 1/n^3. \quad (74)$$

There is a total of ρb_n arriving jobs in a single batch, and for each arriving job

$$\mathbb{P}(\text{the job arrives to queue } i) = \frac{\lambda_i}{\sum_{i=1}^n \lambda_i} \stackrel{(a)}{\leq} \frac{\lambda_i}{(1 - \rho)n} \leq \frac{u_n}{(1 - \rho)n} \stackrel{(b)}{\leq} \frac{1}{2n} \beta_n \leq \frac{1}{2n\hat{\rho}} \beta_n, \quad (75)$$

for all i , where steps (a) and (b) follow from the assumptions that $\sum_{i=1}^n \lambda_i \geq (1 - \rho)n$ (Eq. (10) in Assumption 4.1) and that $u_n \leq \frac{1-\rho}{2}\beta_n$ (in the statement of Theorem 3.4), respectively. From Eq. (75), A_i is stochastically dominated by a binomial random variable $\tilde{A} \stackrel{d}{=} \text{Bino}(\rho b_n, \frac{1}{2n\hat{\rho}}\beta_n)$, with

$$\mathbb{E}(\tilde{A}) = \rho b_n \frac{1}{2n\hat{\rho}}\beta_n = \frac{1}{2} \left(\beta_n \frac{\rho b_n / \hat{\rho}}{n} \right) = \frac{1}{2} \left(\beta_n \frac{m_n}{n} \right) = \frac{1}{2} \hat{u}_n. \quad (76)$$

Based on this expression of $\mathbb{E}(\tilde{A})$, we will now use an exponential tail bound to bound the probability of the event $\{\max_{1 \leq i \leq n} A_i \geq \hat{u}_n\}$. Recall that $b_n = \frac{320}{(1-\rho)^2} \cdot \frac{n \ln n}{\beta_n}$. Using the union bound, we have that

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq i \leq n} A_i \geq \hat{u}_n\right) &= \mathbb{P}(A_i \geq \hat{u}_n, \text{ for some } i) \\ &\leq n\mathbb{P}(A_1 \geq \hat{u}_n) \\ &\leq n\mathbb{P}(\tilde{A} \geq \hat{u}_n) \\ &\stackrel{(a)}{=} n\mathbb{P}(\tilde{A} \geq 2\mathbb{E}(\tilde{A})) \\ &\stackrel{(b)}{\leq} n \exp\left(-\frac{1}{3}\mathbb{E}(\tilde{A})\right) \\ &= n \exp\left(-\frac{\rho}{6\hat{\rho}} \cdot \frac{b_n\beta_n}{n}\right) \\ &\leq n \exp\left(-\frac{\rho}{6} \cdot \frac{b_n\beta_n}{n}\right) \\ &= n \exp\left(-\frac{\rho}{6} \cdot \frac{320}{(1-\rho)^2} \cdot \frac{n \ln n}{\beta_n} \cdot \frac{\beta_n}{n}\right) \\ &\stackrel{(c)}{\leq} n \exp\left(-\frac{160}{6} \ln n\right) \\ &\leq n^{-3}. \end{aligned} \quad (77)$$

Step (a) follows from Eq. (76). Step (b) follows from the following multiplicative form of the Chernoff bound (cf. Chapter 4 of [23]), with $\delta = 1$: $\mathbb{P}(\tilde{A} \geq (1 + \delta)\mu) \leq \exp(-\frac{\delta^2}{2+\delta}\mu)$, where \tilde{A} is a binomial random variable with $\mathbb{E}(\tilde{A}) = \mu$. Step (c) follows from the assumption $\rho \in (1/2, 1)$ (cf. Assumption 4.1), and hence

$$\frac{\rho}{(1-\rho)^2} \geq \rho \geq 1/2. \quad (79)$$

This completes the proof of Lemma 4.5. Q.E.D.

A.8. Proof of Lemma 4.7

Proof. For a set $S \subset \hat{\Gamma}$, denote by $\mathcal{N}^*(S)$ the set of neighbors of S in \hat{G} , i.e., $\mathcal{N}^*(S) = \mathcal{N}(S) \cap \Delta$. To prove Lemma 4.7, we will leverage the fact that the underlying connectivity graph, g_n , is an expander graph with appropriate expansion. As a result, most subsets $S \subset \hat{\Gamma}$ have a large set of neighbors, $\mathcal{N}(S)$, in g_n . Because each server in $\mathcal{N}(S)$ belongs to $\mathcal{N}^*(S)$ independently, as a

consequence of our scheduling policy, we will then use a concentration inequality to show that, with high probability, the sizes of the sets $\mathcal{N}^*(S)$ remain sufficiently large. Using the union bound over the relevant sets S , we will finally conclude that \hat{G} has the desired expansion property, with high probability.

By the definition of a $(\gamma/\hat{u}_n, \hat{u}_n)$ -expander, we are only interested in the expansion of subsets of $\hat{\Gamma}$ with size less than or equal to $|\hat{\Gamma}|\gamma/\hat{u}_n$. We first verify below that the size of such subsets S is sufficiently small to be able to exploit the expansion property of g_n and to infer that $\mathcal{N}^*(S)$ is large. We have

$$\frac{n\gamma/\beta_n}{|\hat{\Gamma}|\gamma/\hat{u}_n} = \frac{n}{|\hat{\Gamma}|} \cdot \frac{\hat{u}_n}{\beta_n} = \frac{n}{m_n} \cdot \frac{\beta_n \frac{m_n}{n}}{\beta_n} = 1, \quad (80)$$

which is equivalent to saying

$$s \leq \gamma n/\beta_n, \quad \forall s \leq |\hat{\Gamma}|\gamma/\hat{u}_n, \quad (81)$$

as desired.

For a set $S \subset \hat{\Gamma}$, we now characterize the size of its neighborhood in \hat{G} , $|\mathcal{N}^*(S)|$, which depends on the distribution of the random subset, Δ . Fix some $s \in \mathbb{N}$ with $s \leq |\hat{\Gamma}|\gamma/\hat{u}_n$. From Eq. (81), we know that $s \leq \gamma n/\beta_n$. Consider some $S \subset \hat{\Gamma}$ with $|S| = s$. Using the expansion property of g_n , we have that $|\mathcal{N}(S)| \geq \beta_n s$. Therefore,

$$\begin{aligned} \mathbb{P}(|\mathcal{N}^*(S)| \leq \hat{u}_n s) &= \mathbb{P}\left(\sum_{j \in \mathcal{N}(S)} \mathbb{I}(j \in \Delta) \leq \hat{u}_n s\right) \\ &\stackrel{(a)}{\leq} \mathbb{P}\left(\text{Bino}\left(|\mathcal{N}(S)|, \frac{b_n}{n}(\rho + 3\epsilon/4)\right) \leq \hat{u}_n s\right) \\ &\stackrel{(b)}{\leq} \mathbb{P}\left(\text{Bino}\left(\beta_n s, \frac{b_n}{n}(\rho + 3\epsilon/4)\right) \leq \hat{u}_n s\right), \end{aligned} \quad (82)$$

for all sufficiently large n . Step (a) follows from the assumption that $\mathbb{P}(j \in \Delta) \geq (\rho + 3\epsilon/4)\frac{b_n}{n}$, and step (b) from the inequality $|\mathcal{N}(S)| \geq \beta_n s$. We observe that

$$\begin{aligned} \mu &\triangleq \mathbb{E}\left(\text{Bino}\left(\beta_n s, \frac{b_n}{n}(\rho + 3\epsilon/4)\right)\right) \\ &= (\rho + 3\epsilon/4) \frac{\beta_n b_n}{n} s \\ &\stackrel{(a)}{=} (\rho + 3\epsilon/4) \frac{1}{n} \cdot \frac{80}{\epsilon^2} \cdot \frac{n \ln n}{\beta_n} \beta_n s \\ &= (\rho + 3\epsilon/4) \frac{80 \ln n}{\epsilon^2} s, \end{aligned} \quad (83)$$

where in step (a) we used the substitution $b_n = \frac{80}{\epsilon^2} \cdot \frac{n \ln n}{\beta_n}$. We also have that

$$\hat{u}_n = \beta_n \frac{m_n}{n}$$

$$\begin{aligned}
&= \beta_n \frac{\rho b_n}{\hat{\rho} n} \\
&= \beta_n \frac{\rho}{\hat{\rho} n} \cdot \frac{80}{\epsilon^2} \cdot \frac{n \ln n}{\beta_n} \\
&= \frac{\rho}{\hat{\rho}} \cdot \frac{80 \ln n}{\epsilon^2}.
\end{aligned} \tag{84}$$

By combining Eqs. (83) and (84), we can derive a useful lower bound on the quantity $1 - \frac{s\hat{u}_n}{\mu}$, which is recorded in the lemma that follows.

Lemma A.4 *We have that*

$$1 - \frac{s\hat{u}_n}{\mu} \geq \frac{\epsilon}{2}. \tag{85}$$

Proof. Using Eqs. (83) and (84) in the first step below, we have that

$$1 - \frac{s\hat{u}_n}{\mu} = 1 - \frac{\rho}{\hat{\rho}(\rho + 3\epsilon/4)}.$$

Recall that $\epsilon = (1 - \rho)/2$, so that $\rho = 1 - 2\epsilon$ and that $\hat{\rho} = 1/(1 + \epsilon/4)$. Using these substitutions, we obtain

$$\begin{aligned}
1 - \frac{s\hat{u}_n}{\mu} &= 1 - \frac{(1 - 2\epsilon)(1 + \epsilon/4)}{1 - 2\epsilon + 3\epsilon/4} \\
&= \frac{3\epsilon/4 - \epsilon/4 + 2\epsilon^2/4}{1 - 5\epsilon/4} \\
&= \frac{\epsilon(1 + \epsilon)/2}{1 - 5\epsilon/4} \\
&\geq \frac{\epsilon}{2}.
\end{aligned}$$

Q.E.D.

To obtain an upper bound for the probability in Eq. (82), we substitute Eqs. (83) and (85) into Eq. (82). Given the assumption that $s \leq \gamma n / \beta_n$, we have that

$$\begin{aligned}
\mathbb{P}(|\mathcal{N}^*(S)| \leq \hat{u}_n s) &\leq \mathbb{P}\left(\text{Bino}\left(\beta_n s, \frac{b_n}{n}(\rho + 3\epsilon/4)\right) \leq \hat{u}_n s\right) \\
&\stackrel{(a)}{\leq} \exp\left(-\frac{1}{2} \left(\frac{\epsilon}{2}\right)^2 \mu\right) \\
&\stackrel{(b)}{=} \exp\left(-\frac{\epsilon^2}{8} \cdot \frac{80 \ln n}{\epsilon^2} (\rho + 3\epsilon/4) s\right) \\
&= \exp(-(10 \ln n)(\rho + 3\epsilon/4) s) \\
&\stackrel{(c)}{\leq} \exp(-5 \ln n) s \\
&= \frac{1}{n^{5s}}.
\end{aligned} \tag{86}$$

for all sufficiently large n . Step (a) is based on a multiplicative form of the Chernoff bound (cf. Chapter 4 of [23]), $\mathbb{P}(X \leq (1 - \delta)\mu) \leq \exp(-\frac{1}{2}\delta^2\mu)$, where X is a binomial random variable with $\mathbb{E}(X) = \mu$, and

$$\delta = 1 - \frac{s\hat{u}_n}{\mu} \geq \epsilon/2, \quad (87)$$

where the last inequality follows from Lemma A.4. Step (b) follows from Eq. (83), and (c) from the assumption that $\rho \geq 1/2$.

We now apply Eq. (86) to subsets of $\hat{\Gamma}$, and use the union bound. We have, for all sufficiently large n , that

$$\begin{aligned} \mathbb{P}\left(\hat{G} \text{ is not a } (\gamma/\hat{u}_n, \hat{u}_n)\text{-expander}\right) &\leq \mathbb{P}(\exists S \subset \hat{\Gamma} \text{ such that: } |S| \leq |\hat{\Gamma}|\gamma/\hat{u}_n \text{ and } |\mathcal{N}^*(S)| \leq \hat{u}_n|S|) \\ &\stackrel{(a)}{\leq} \sum_{s=1}^{|\hat{\Gamma}|\gamma/\hat{u}_n} \left(\sum_{S \subset \hat{\Gamma}, |S|=s} \mathbb{P}(|\mathcal{N}^*(S)| \leq \hat{u}_n s) \right) \\ &\leq \sum_{s=1}^{|\hat{\Gamma}|\gamma/\hat{u}_n} \binom{|\hat{\Gamma}|}{s} \mathbb{P}(|\mathcal{N}^*(S)| \leq \hat{u}_n s) \\ &\stackrel{(b)}{<} \sum_{s=1}^{|\hat{\Gamma}|\gamma/\hat{u}_n} b_n^s \mathbb{P}(|\mathcal{N}^*(S)| \leq \hat{u}_n s) \\ &\stackrel{(c)}{\leq} \sum_{s=1}^{|\hat{\Gamma}|\gamma/\hat{u}_n} b_n^s \frac{1}{n^{5s}} \\ &\leq \sum_{s=1}^{\infty} (b_n/n^5)^s \\ &= \frac{b_n/n^5}{1 - b_n/n^5}. \end{aligned} \quad (88)$$

Step (a) is the union bound. In step (b), we used the bound $\binom{n}{k} \leq n^k$, and the fact that $|\hat{\Gamma}| = m_n = \frac{\rho}{\beta} b_n < b_n$. Step (c) follows from Eq. (86). Because $\beta_n \gg \ln n$, we have that $b_n \lesssim \frac{n \ln n}{\beta_n} \ll n$, and hence

$$\frac{b_n}{n^5} \leq \frac{1}{n^3}, \quad (89)$$

for all sufficiently large n . Combining Eqs. (88) and (89), we conclude that

$$\mathbb{P}\left(\hat{G} \text{ is not a } \left(\frac{\gamma}{\hat{u}_n}, \hat{u}_n\right)\text{-expander}\right) \leq \frac{1}{n^3}, \quad (90)$$

for all sufficiently large n . This proves our claim. Q.E.D.

Appendix B: Expanded Modular Architectures

In this appendix, we start by describing the graph product, and subsequently we discuss the implications of using an expander graph.

Construction of the Architecture. We first express the average degree as a product, $d_n = d_n^m \cdot d_n^e$, where the relative magnitudes of d_n^m and d_n^e are a design choice. The architecture is constructed as follows.

1. Similar to the case of the Modular architecture, partition I and J into equal-sized clusters of size d_n^m . We will refer to the index set of the queue and server clusters as \mathcal{Q} and \mathcal{S} , respectively. For any $i \in I$ and $j \in J$, denote by $q(i) \in \mathcal{Q}$ and $s(j) \in \mathcal{S}$, the indices of the queue and server clusters to which i and j belong, respectively.
2. Let g_n^e be a bipartite graph of maximum degree d_n^e whose left and right nodes are the queue and server clusters, \mathcal{Q} and \mathcal{S} , respectively. Let E^e be the set of edges of g_n^e .
3. To construct the interconnection topology $g_n = (I \cup J, E)$, let $(i, j) \in E$ if and only if their corresponding queue and server clusters are connected in g_n^e , i.e., if $(q(i), s(j)) \in E^e$.

Note that by the above construction, each queue is connected to at most d_n^e server clusters through g_n^e , and within each connected cluster, to d_n^m servers. Therefore, the maximum degree of g_n is $d_n^m \cdot d_n^e = d_n$.

Scheduling Policy. The scheduling policy requires the knowledge of the arrival rate vector, λ_n , and involves two stages. For a given λ_n , the computation in the first stage is performed only once, while the steps in the second stage are repeated throughout the operation of the system.

1. Compute a feasible flow, $\{f_{q,s}\}_{(q,s) \in E^e}$, over the graph g_n^e , where the incoming flow at each queue cluster $q \in \mathcal{Q}$ is equal to $\sum_{i \in q} \lambda_i$, and the outgoing flow at each server cluster $s \in \mathcal{S}$ is constrained to be less than or equal to $\frac{1+\rho}{2} d_n^m$. (It turns out that, under our assumptions, such a feasible flow exists [33].) Denote by $f_{q,s}$ the total rate of flow from the queue cluster q to the server cluster s .
2. Arriving jobs first wait in queue until they are fetched by a server. When a server becomes available, it chooses a neighboring queue cluster (w.r.t. the topology of g_n^e) with probability roughly proportional to the flow between the clusters. In particular, a server in cluster s chooses the queue cluster q with probability

$$p_{s,q} = \frac{f_{q,s}}{\sum_{q' \in \mathcal{N}(s)} f_{q',s}} \cdot \frac{1+\rho}{2} + \frac{1}{\deg(s)} \cdot \frac{1-\rho}{2}, \quad (91)$$

where $\deg(s)$ is the degree of s in g_n^e . Within the chosen cluster, the server starts serving a job from an arbitrary non-empty queue, or, if all queues in the cluster are empty, the server initiates an idling period whose length is exponentially distributed with mean 1.

When the graph g_n^e is an expander graph, we refer to the topology created via the above procedure as an *Expanded Modular architecture generated by g_n^e* .

Note that an Expanded Modular architecture is constructed as a “product” between an expander graph across the queue and server clusters, and a fully connected graph for each pair of connected

clusters. As a result, its performance is also of a hybrid nature: the expansion properties of g_n^e guarantee a large capacity region, while a diminishing delay is obtained as a result of the growing size of the server and queue clusters. We summarize this in the next theorem. Here we assume that d_n^e is sufficiently large so that the expander graph described in Lemma 3.3 exists. The reader is referred to Section 3.4.5 of [33] for the proof of the theorem (although with different choices for some of the constants).

Theorem B.1 (Capacity and Delay of Expanded Modular Architectures) *Suppose that $d_n = d_n^m \cdot d_n^e$. Let $\gamma = \sqrt{\rho}$ and $\beta_n = \frac{1}{2} \cdot \frac{\ln(1/\rho)}{1 + \ln(1/\rho)} d_n^e$. Let g_n^e be a $(\gamma/\beta_n, \beta_n)$ -expander with maximum degree d_n^e , and let g_n be an Expanded Modular architecture generated by g_n^e . If*

$$u_n \leq \frac{1 + \rho}{2} \beta_n = \frac{1 + \rho}{4} \cdot \frac{\ln(1/\rho)}{1 + \ln(1/\rho)} d_n^e, \quad (92)$$

then, under the scheduling policy described above, we have that

$$\sup_{\lambda_n \in \Lambda_n(u_n)} \mathbb{E}(W | \lambda_n) \lesssim \frac{c}{d_n^m}, \quad (93)$$

where c is a constant that does not depend on n .

A Tradeoff between the Size of the Capacity Region and the Delay. For the Expanded Modular architecture, the relative values of d_n^m and d_n^e reflect a design choice: a larger value of d_n^e ensures a larger capacity region, while a larger value of d_n^m yields smaller delays. Therefore, while the Expanded Modular architecture is able to provide a strong delay guarantee that applies to *all* arrival rate vectors in $\Lambda_n(u_n)$, it comes at the expense of either a slower rate of diminishing delay (small d_n^m) or a smaller capacity region (small d_n^e).