



Resource Allocation in One-dimensional Distributed Service Networks

Nitish K Panigrahy, Prithwish Basu, Philippe Nain, Don Towsley, Ananthram Swami, Kevin S. Chan, Kin K Leung

► To cite this version:

Nitish K Panigrahy, Prithwish Basu, Philippe Nain, Don Towsley, Ananthram Swami, et al.. Resource Allocation in One-dimensional Distributed Service Networks. MASCOTS 2019 - 27th IEEE International Symposium on the Modeling, Analysis, and Simulation of Computer and Telecommunication Systems, Oct 2019, Rennes, France. pp.14-26, 10.1109/MASCOTS.2019.00013 . hal-02267631v2

HAL Id: hal-02267631

<https://inria.hal.science/hal-02267631v2>

Submitted on 17 Feb 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Resource Allocation in One-dimensional Distributed Service Networks

Nitish K. Panigrahy[†], Prithwish Basu[¶], Philippe Nain[¶], Don Towsley[†], Ananthram Swami[‡],
Kevin S. Chan[‡] and Kin K. Leung[§]

[†] University of Massachusetts Amherst, MA, USA. Email: {nitish, towsley}@cs.umass.edu

[¶] Raytheon BBN Technologies, Cambridge, MA 02138, USA. Email: prithwish.basu@raytheon.com

[¶] Inria, 06902 Sophia Antipolis Cedex, France. Email: philippe.nain@inria.fr

[‡] Army Research Laboratory, Adelphi, MD 20783, USA. Email: {ananthram.swami, kevin.s.chan}.civ@mail.mil

[§] Imperial College London, London SW72AZ, UK. Email: kin.leung@imperial.ac.uk

Abstract—We consider assignment policies that allocate resources to users, where both resources and users are located on a one-dimensional line $[0, \infty)$. First, we consider unidirectional assignment policies that allocate resources only to users located to their left. We propose the Move to Right (*MTR*) policy, which scans from left to right assigning the nearest available resource located to the right of a user, and contrast it to the Unidirectional Gale-Shapley (*UGS*) matching policy. While both policies among all unidirectional policies, minimize the expected distance traveled by a request, *MTR* is fairer. Moreover, we show that when user and resource locations are modeled by statistical point processes, and resources are allowed to satisfy more than one user, the spatial system under unidirectional policies can be mapped into bulk service queueing systems, thus allowing the application of many queueing theory results that yield closed form expressions. As we consider a case where different resources can satisfy different numbers of users, we also generate new results for bulk service queues. We also consider bidirectional policies where there are no directional restrictions on resource allocation and develop an algorithm for computing the optimal assignment which is more efficient than known algorithms in the literature when there are more resources than users. Finally, numerical evaluation of performance of unidirectional and bidirectional allocation schemes yields design guidelines beneficial for resource placement.

I. INTRODUCTION

The past few years have witnessed significant growth in the use of distributed network analytics involving agile code, data and computational resources. In many such networked systems, for example, Internet of Things [4], a large number of computational and storage resources are widely distributed in the physical world. These resources are accessed by various end users/applications that are also distributed over the physical space. Assigning users or applications to resources efficiently is key to sustained high-performance operation of the system.

In some systems, requests are transferred over a communication network to a server that provides a needed resource. In other systems, servers are mobile and physically move to the user making a request. Examples of the former type of service include accessing storage resources over a wireless network to store files and requesting computational resources to run image processing tasks; an example of the latter type of service is

the arrival of ride-sharing vehicles to the user's location over a road transportation network.

Not surprisingly, the spatial distribution of resources and users¹ in the network is an important factor in determining the overall performance of the service. A key measure of performance is *expected request distance*, that is average of the pairwise distances between each user and its allocated resource/server (where distance is measured on the network). This directly translates to latency incurred by a user when accessing the service, which is arguably among the most important criteria in distributed service applications. For example, in wireless networks, signal attenuation is strongly coupled to request distance, therefore developing allocation policies to minimize request distance can help reduce energy consumption, an important concern in battery-operated wireless networks. Another important practical constraint in distributed service networks is *service capacity*. For example, in network analytics applications, a networked storage device can only support a finite number of concurrent users; similarly, a computational resource can only support a finite number of concurrent processing tasks. Likewise, in physical service applications like ride-sharing, a vehicle can pick up a finite number of passengers at once.

Therefore, a primary problem in such distributed service networks is to efficiently assign each user to a suitable resource so as to minimize expected request distance and ensure no resource is allocated to more users than its capacity. If the entire system is being managed by a single administrative entity such as a ride sharing service, or a datacenter network where analytics tasks are being assigned to available CPUs, there are economic benefits in minimizing the expected request distance across all (user, resource) pairs, which is tantamount to minimizing the average delay in the system.

The general version of this capacitated assignment problem can be solved by modeling it as a *minimum cost flow* problem on graphs [3] and running the *network simplex algorithm* [17]. However, if the network has a low-dimensional structure and

¹We use the terms “users” and “requesters” interchangeably and same holds true for the terms “resources” and “servers”.

some assumptions about the spatial distributions of users and resources hold, more efficient methods can be developed [12].

In this paper, we consider two one-dimensional network scenarios that motivate the study of this special case of the user-to-resource assignment problem.

The first scenario is ride-hailing on a one-way street where vehicles move right to left. If the vehicles of a ride-sharing company are distributed along the street at a certain time, and users equipped with smartphone ride-hailing apps request service, the system attempts to assign vehicles with spare capacity located towards the right of the users so as to minimize expected “pick up” distance. Abadi et al. [1] introduced this problem and presented a policy known as Unidirectional Gale-Shapley² matching (*UGS*) to minimize expected pick up distance. In this policy, all users concurrently emit rays of light toward their right and each user is matched with the vehicle that first receives its emitted ray. While the well-known Gale-Shapley matching algorithm [7] matches user-resource pairs that are mutually nearest to each other, its unidirectional variant, *UGS*, matches a user to the nearest resource on its right. Note that, this one-dimensional network setting also applies to vehicular wireless ad-hoc networks on a one-lane roadway [10], [15]³, where users are in vehicles and servers are attached to fixed infrastructure such as lamp posts. Users attempt to allocate their computation tasks over the wireless network to servers located to their right so that they can retrieve the results with little effort while driving by.

In this paper, we propose another policy “Move to Right” policy (or *MTR*) which has the same “expected distance traveled by a request” (*request distance*) as *UGS* but has a lower variance. *MTR* sequentially allocates users to the geographically nearest available vehicle located to his/her right. When user and resource locations are modeled by statistical point processes, the one-dimensional unidirectional space behaves similar to time and notions from queueing theory can be applied. In particular, when user and vehicle locations are modeled by independent Poisson processes, expected request distance can be characterized in closed form by considering inter-user and inter-server distances as parameters of a *bulk service M/M/1* queue where the bulk service capacity denotes the maximum number of users that can be handled by a server. We equate request distance in the spatial system to the expected *sojourn time* in the corresponding queueing model⁴. This natural mapping allows us to use well-known results from queueing theory and in some cases to propose new queueing theoretic models to characterize request distances for a number of interesting situations beyond *M/M/1* queues.

The second scenario involves a convoy of vehicles traveling on a one-dimensional space, for example, trucks on a highway or boats on a river. Some vehicles have expensive camera sensors (image/video) but have inadequate computational stor-

age or processing power. On the other hand, cheap storage and processing is easily available on several other vehicles. The cameras periodically take photos/videos as they move through space and want them processed / stored. In such cases, bidirectional assignment schemes are more suitable. Since no directionality restrictions are imposed on the allocation algorithms, computing the optimal assignment is not as simple as in the unidirectional case.

We explore the special structure of the one-dimensional topology to develop an optimal algorithm that assigns a set of requesters R to a set of resources S such that the total assignment cost (and hence the expected request distance) is minimized. This problem has been recently solved for $|R| = |S|$ [6] with time complexity $O(|R|)$. Note that other assignment algorithms in literature such as the Hungarian primal-dual algorithm and Agarwal’s variant [2] have time complexities $O(|R|^3)$ and $O(|R|^{2+\epsilon})$ respectively and assume $|R| = |S|$ for general and Euclidean distance measures. However, we are interested in the case when $|R| < |S|$. We propose a Dynamic Programming based algorithm which solves this case with time complexity $O(|R|(|S| - |R| + 1))$.

Our contributions are summarized below:

- 1) Analysis of simple unidirectional allocation policies *MTR* and *UGS* yielding closed form expressions for expected request distance.
 - When inter-requester and inter-resource distances are exponentially distributed, we model unidirectional policies as a bulk service *M/M/1* queue.
 - When inter-requester distances are generally distributed but the inter-resource distances are exponentially distributed, we model the situation using an accessible batch service *G/M/1* queue.
 - When inter-requester distances are exponentially distributed but inter-resource distances are generally distributed, we model the spatial system as an accessible batch service *M/G/1* queue with the first batch having exceptional service time. To the best of our knowledge this system has not been studied previously in the queueing theory literature.
 - We include several generalizations of our framework. In the first place we discuss a simulation driven conjecture for evaluating request distance for general distance distributions under heavy traffic. We also investigate the heterogeneous server capacity scenario where server capacity is a random variable and to the best of our knowledge this system has not been studied previously in the queueing theory literature. We derive expressions for expected request distance when servers have infinite capacity.
- 2) A novel algorithm for optimal (bidirectional) assignment with time complexity $O(|R|(|S| - |R| + 1))$.
- 3) A numerical and simulation study of different assignment policies: *UGS*, *MTR*, a bi-directional heuristic allocation policy (Gale-Shapley) and the optimal policy.

The paper is organized as follows. The next section dis-

²We rename *queue matching* defined in [1] as Unidirectional Gale-Shapley Matching to avoid overloading the term *queue*.

³Furthermore, [10] confirms that vehicle location distribution on the streets in Central London can be closely approximated by a Poisson distribution.

⁴Sojourn time is the sum of waiting and service times in a queue.

cusses related work. Section III contains technical preliminaries. We show the equivalence of UGS and MTR w.r.t expected request distance in Section IV, and present results associated with the case when servers are Poisson distributed in Section V. In Section VI, we develop formulations for expected request distance when either user or server placements are described by Poisson processes. We include some generalizations of our framework such as analysis under general distance distributions, results for heterogeneous server capacity and uncapacitated allocation in Section VII. The optimal bidirectional allocation strategy is presented in Section VIII. We compare the performance of various local allocation strategies in Section IX. We conclude the paper in Section X.

II. RELATED WORK

Poisson Matching: Holroyd et al. [11] first studied translation invariant matchings between two d -dimensional Poisson processes with equal densities. Their primary focus was obtaining upper and lower bounds on expected matching distance for stable matchings. Abadi et al. [1] introduced “Unidirectional Gale-Shapley” matching (UGS) and derived bounds on the expected matching distance for stable matchings between two one-dimensional Poisson processes with different densities. In this paper, we propose another unidirectional allocation policy: “Move To Right” policy (MTR) and provide explicit expressions for the expected matching distance for both MTR and UGS when either requesters or servers are distributed according to a renewal process and the other according to a Poisson process.

Exceptional Queueing Systems and Accessible Batches: Welch et al. [20] first studied an M/G/1 queue where a customer arriving when the server is idle has a different service time than the others. Bulk service M/G/1 queues has been studied in [5]. Authors in [8] analyzed a bulk service G/M/1 queue with accessible or non-accessible batches where an accessible batch is a batch in service allowing subsequent arrivals, while the service is on. In this work, we model the spatial system using an accessible batch service queue with the first batch having exceptional service time. To the best of our knowledge this system has not been studied previously in queueing theory literature.

Euclidean Bipartite Matching: The optimal user-server assignment problem can be modeled as a minimum-weight matching on a weighted bipartite graph where weights on edges are given by the Euclidean distances between the corresponding vertices [16]. Well-known polynomial time solutions exist for this problem, such as the modified Hungarian algorithm proposed by Agarwal et al. [2] with a running time of $O(|R|^{2+\epsilon})$, where $|R|$ is the total number of users. In the case of an equal number of users and servers, the optimal user-server assignment on a real line is known [6]. In this paper, we consider the case when there are fewer users than servers.

III. TECHNICAL PRELIMINARIES

Consider a set of users R and a set of servers S . Each user makes a request that can be satisfied by any server. Assume

that each server $j \in S$ has capacity $c_j \in \mathbb{Z}^+$ corresponding to the maximum number of requests that it can process. Suppose users and servers are located on a line \mathcal{L} . Formally, let $r : R \rightarrow \mathcal{L}$ and $s : S \rightarrow \mathcal{L}$ be the location functions for users and servers, respectively, such that a distance $d_{\mathcal{L}}(r, s)$ is well defined for all pairs $(r, s) \in R \times S$. Initially we assume that all servers have equal capacities i.e. $c_j = c \forall j \in S$. Later in Section VII-B we extend our analysis to a case in which server capacities are integer random variables.

A. User and server spatial distributions

Let $0 \leq r_1 \leq r_2 \leq \dots$ represent user locations and $0 \leq s_1 \leq s_2 \leq \dots$ be the server locations. Let $X_j = s_j - s_{j-1}, j \geq 1, s_0 = 0$, denote the inter-server distances and $Y_i = r_i - r_{i-1}, i \geq 1, r_0 = 0$, the inter-user distances. We assume $\{X_j\}_{j \geq 1}$ to be a renewal process with cumulative distribution function (cdf)

$$\mathbb{P}(X_j \leq x) = F_X(x). \quad (1)$$

We also assume $\{Y_i\}_{i \geq 1}$ to be a renewal process with cdf $F_Y(x)$, i.e.,

$$\mathbb{P}(Y_i \leq x) = F_Y(x). \quad (2)$$

We denote $\alpha_X = 1/\mu$ and σ_X^2 to be the mean and variance associated with F_X . Similarly let $\alpha_Y = 1/\lambda$ and σ_Y^2 be the mean and variance associated with F_Y . We let $\rho = \lambda/\mu$ and assume that $\rho < c$. Denote by $F_X^*(s) = \int_0^\infty e^{-sx} dF_X(x)$ and $F_Y^*(s)$ the Laplace-Stieltjes transform (LST) of F_X and F_Y with $s \geq 0$.

In our paper, we consider various inter-server and inter-user distance distributions, including exponential, deterministic, uniform and hyperexponential.

B. Allocation policies

One of our goals is to analyze the performance of various request allocation policies using expected request distance as a performance metric. We define various allocation policies as follows.

- **Unidirectional Gale-Shapley (UGS):** In UGS, each user simultaneously emits a ray to their right. Once the ray hits an unallocated server s , the user is allocated to s .
- **Move To Right (MTR):** In MTR, starting from the left, each user is allocated sequentially to the nearest available server to its right.
- **Gale-Shapley (GS) [7]:** In this matching, each user selects the nearest server and each server selects its nearest user. Remove reciprocating pairs, and continue.
- **Optimal Matching:** This matching minimizes average request distance among all feasible allocation policies.

IV. UNIDIRECTIONAL ALLOCATION POLICIES

In this Section, we establish the equivalence of UGS and MTR w.r.t number of requests that traverse a point and expected request distance. Define N_x^P and D_i^P to be random variables for the number of requests that traverse point $x \in \mathcal{L}$ and distance between user i and its allocated server under

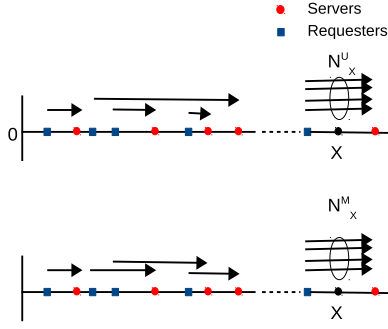


Fig. 1: Allocation of users to servers on the one-dimensional network. Top: UGS, Bottom: MTR allocation policy.

policy P , respectively. Thus N_x^U and N_x^M denote the number of requests that traverse point $x \in \mathcal{L}$ under UGS and MTR, respectively, as shown in Figure 1. Consider the following definition of busy cycle in a service network.

Definition 1. A busy cycle for a policy P is an interval $I = [a, b] \subset \mathcal{L}$ such that $\exists i, j$ with $r_i = a, s_j = b$ for which $N_x^P > 0, \forall x \in I$ and $N_x^P = 0$ for $x = a - \epsilon$ and $x = b + \epsilon$ with ϵ being an infinitesimal positive value.

We have the following theorem.

Theorem 1. $N_x^U = N_x^M, x \geq 0$.

Proof. See [18]. \square

Corollary 1. $\mathbb{E}[D^U] = \mathbb{E}[D^M]$ i.e. the expected request distances are the same for both UGS and MTR under steady state.

Proof. Under steady state both N_x^U and N_x^M converge to a random variable. Applying Little's law we have $\mathbb{E}[D^U] = \mathbb{E}[D^M]$. \square

Remark 1. Note that Theorem 1 applies to any inter-server or inter-user distance distribution. It also applies to the case where servers have capacity $c > 1$.

Remark 2. Although MTR and UGS are equivalent w.r.t. the expected request distance, MTR tends to be fairer, i.e., has low variance⁵ for expected request distance.

V. UNIDIRECTIONAL POISSON MATCHING

In this section, we characterize request distance statistics under unidirectional policies when both users and servers are distributed according to two independent Poisson processes. We first analyze MTR as follows.

A. MTR

Under this allocation policy, the service network can be modeled as a bulk service M/M/1 queue. A bulk service M/M/1 queue provides service to a group of c or fewer customers. The server serves a bulk of at most c customers

⁵It is well known in queueing theory that among all service disciplines the variance of the waiting time is minimized under FCFS policy [13]. In Section V we show that MTR maps to a temporal FCFS queue.

whenever it becomes free. Also customers can join an existing service if there is room which is an example of accessible batch. In Section VI we describe the notion of accessible batches in greater detail. The service time for the group is exponentially distributed and customer arrivals are described by a Poisson process. The distance between two consecutive users in the service network can be thought of as inter-arrival time between customers in the bulk service M/M/1 queue. The distance between two consecutive servers maps to a bulk service time.

Having established an analogy between the service network and the bulk service M/M/1 queue, we now define the state space for the service network. Consider the definition of N_x as the number of requests⁶ that traverse point $X \in \mathcal{L}$ under MTR. In steady state, N_x converges to a random variable N provided $\lambda < c\mu$. Let π_k denote $\mathbb{P}[N = k]$ with $k \geq 0$.

Following the procedure in [14], we obtain the steady state probability vector $\pi = [\pi_i, i \geq 0]$. In the service network, request distance corresponds to the sojourn time in the bulk service M/M/1 queue. By applying Little's formula, we obtain the following expression for the expected request distance

$$\mathbb{E}[D] = \frac{r_0}{\lambda(1 - r_0)}, \quad (3)$$

where r_0 is the only root in the interval $(0, 1)$ of the following equation (with r as the variable)

$$\mu r^{c+1} - (\lambda + \mu)r + \lambda = 0. \quad (4)$$

1) When server capacity: $c = 1$: When $c = 1$, $r_0 = \rho$ is a solution of (4). Thus we can evaluate the expected request distance as

$$\mathbb{E}[D] = \frac{\rho}{\lambda(1 - \rho)} = \frac{1}{\mu - \lambda}. \quad (5)$$

Note that, when server capacity is one, the service network can be modeled as an M/M/1 queue. In such a case, (5) is the mean sojourn time for an M/M/1 queue.

B. UGS

When both users and servers are Poisson distributed and servers have unit capacity, the request distance in UGS has the same distribution as the busy cycle in the corresponding Last-Come-First-Served Preemptive-Resume (LCFS-PR) queue having the density function [1]

$$f_{D^U}(x) = \frac{1}{x\sqrt{\rho}} e^{(\lambda+\mu)x} I_1(2x\sqrt{\lambda\mu}), \quad x > 0, \quad (6)$$

where $\rho = \lambda/\mu$ and I_1 is the modified Bessel function of the first kind. Thus the expected request distance is equivalent to the average busy cycle duration in a LCFS-PR queue given by $1/(\mu - \lambda)$ [1].

When servers have capacities $c > 1$ it is difficult to characterize the expected request distance explicitly. However, by Theorem 1, the expected request distance under UGS is the same as that of MTR given by (3).

Distribution	Parameters	$F_X(x)$	$B(x)$
Exponential	μ : rate	$1 - e^{-\mu x}$	$\frac{1}{\lambda} [1 - e^{-\lambda x}] - \frac{1}{\lambda + \mu} [1 - e^{-(\lambda + \mu)x}]$
Uniform	b : maximum value	$x/b, 0 \leq x \leq b$	$\frac{1}{\lambda^2 b} [1 - e^{-\lambda b}] - \frac{e^{-\lambda x}}{\lambda}$
Deterministic	d_0 : constant	$1, x \geq d_0$	$\frac{e^{-\lambda d_0} - e^{-\lambda x}}{\lambda}$
Hyper -exponential	l : order p_j : phase probability μ_j : phase rate	$1 - \sum_{j=1}^l p_j e^{-\mu_j x}$	$\frac{1}{\lambda} [1 - e^{-\lambda x}] - \sum_{j=1}^l \frac{p_j}{\lambda + \mu_j} [1 - e^{-(\lambda + \mu_j)x}]$

TABLE I: Properties of specific inter-server distance distributions.

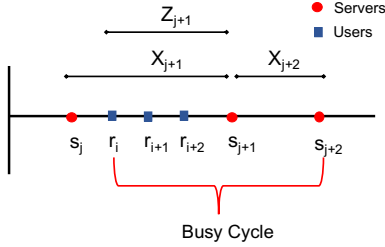


Fig. 2: Allocation of users to servers under MTR policy.

VI. UNIDIRECTIONAL GENERAL MATCHING

We now derive expressions for the expected request distance when either users or servers are distributed according to a Poisson process and the other by renewal process.

A. Notion of exceptional service and accessible batches

We discuss the notion of exceptional service and accessible batches applicable to our service network as follows. Consider a service network with $c = 2$ as shown in Figure 2. Consider a user r_i . Let s_j be the server immediately to the left of r_i . We assume all users prior to r_i have already been allocated to servers $\{s_k, 1 \leq k \leq j\}$. MTR allocates both r_i and r_{i+1} to s_{j+1} and allocates r_{i+2} to s_{j+2} . We denote $[r_i, s_{j+2}]$ as a busy cycle of the service network. We have the following queueing theory analogy.

User r_i can be thought of as the first customer in a queueing system that initiates a busy period while r_{i+1} sees the system busy when it arrives. Because only r_i is in service at the arrival of r_{i+1} , r_{i+1} enters service with r_i and the two customers form a batch of size 2. and depart at time s_{j+1} . This is an example of an *accessible batch* [8]. Recall that an accessible batch admits subsequent arrivals, while the service is on, until the server capacity c is reached.

The service time for the batch, r_i, r_{i+1} , is described by the random variable Z_{j+1} which is different or *exceptional* when compared to service times of successive batches such as the one consisting of r_{i+2} . The service time for the second batch is X_{j+2} . Note that, Z_{j+1} only depends on X_{j+2} and Y_{i+2} . Thus when either X_{j+2} or Y_{i+2} is described by a Poisson process and the other by renewal process, Z_{j+1} converges to a random variable Z under steady state conditions. Denote $F_Z(x)$ and $f_Z(x)$ as the distribution and density functions

for the random variable Z . Thus the service network can be mapped to an exceptional service with accessible batches queueing (ESABQ) model. We formally define ESABQ as follows.

ESABQ: Consider a queueing system where customers are served in batches of maximum size c . A customer entering the queue and finding fewer than c customers in the system joins the current batch and enters service at once, otherwise it joins a queue. After a batch departs leaving k customers in the buffer, $\min(c, k)$ customers form a batch and enter service immediately. There are two different service times cdfs, $F_Z(x)$ (exceptional batch) with mean $\alpha_Z = 1/\mu_Z$ and $F_X(x)$ (ordinary batch) with mean $\alpha_X = 1/\mu$. A batch is exceptional if its oldest customer entered an empty system, otherwise it is a regular batch. When the service time expires, all customers in the server depart at once, regardless of the nature of the batch (exceptional or regular).

1) *Evaluation of the distribution function:* $F_Z(x)$: In this Section, we compute explicit expressions for the distribution function $F_Z(x)$ applicable to our service network.

When $F_X(x) \sim \text{Expo}(\mu)$: In this case, we invoke the memoryless property of the exponential distribution F_X . Thus the exceptional distribution, F_Z , is

$$F_Z(x) = F_X(x) = 1 - e^{-\mu x}, x \geq 0. \quad (7)$$

When $F_Y(x) \sim \text{Expo}(\lambda)$: Using the memoryless property of F_Y , F_Z can be computed as

$$\begin{aligned} F_Z(x) &= \Pr(X - Y < x | Y < X) = \Pr(X - Y < x | X - Y > 0) \\ &= \frac{\Pr(X - Y < x) - \Pr(X - Y < 0)}{1 - \Pr(X - Y < 0)} \\ &= \frac{D_{XY}(x) - D_{XY}(0)}{1 - D_{XY}(0)}, x \geq 0, \end{aligned} \quad (8)$$

where $D_{XY}(x)$ is the distribution of the random variable $X - Y$ (also known as difference distribution). $D_{XY}(x)$ can be expressed as

$$\begin{aligned} D_{XY}(x) &= \Pr(X - Y \leq x) = \int_0^\infty \Pr(X - y \leq x) \Pr(Y = y) dy \\ &= \int_0^\infty F_X(x + y) \lambda e^{-\lambda y} dy = \int_x^\infty F_X(z) \lambda e^{-\lambda(z-x)} dz \\ &= \lambda e^{\lambda x} \left[\int_0^\infty F_X(z) e^{-\lambda z} dz - \int_0^x F_X(z) e^{-\lambda z} dz \right] \\ &= \lambda e^{\lambda x} [A(F_X) - B(x)], \end{aligned} \quad (9)$$

⁶We drop the superscript (M) for brevity.

where \mathcal{A} is the Laplace Transform operator on the function F_X and $\mathcal{B}(x)$ is denoted by

$$\mathcal{B}(x) = \int_0^x F_X(z) e^{-\lambda z} dz \quad (10)$$

Clearly $\mathcal{B}(0) = 0$. Thus combining (8) and (9) yields

$$F_Z(x) = \frac{\lambda e^{\lambda x} [\mathcal{A}(F_X) - \mathcal{B}(x)] - \lambda \mathcal{A}(F_X)}{1 - \lambda \mathcal{A}(F_X)}, \quad (11)$$

$$f_Z(x) = \frac{\lambda^2 e^{\lambda x} [\mathcal{A}(F_X) - \mathcal{B}(x)] - \lambda F_X(x)}{1 - \lambda \mathcal{A}(F_X)}, \quad (12)$$

$$\alpha_Z = \int_0^\infty x f_Z(x) dx, \quad \sigma_Z^2 = \left[\int_0^\infty x^2 f_Z(x) dx \right] - \alpha_Z^2. \quad (13)$$

Expressions for $\mathcal{B}(x)$ are presented in Table I. We can evaluate $\mathcal{A}(F_X)$ by setting $\mathcal{A}(F_X) = \mathcal{B}(\infty)$. Detailed derivations are relegated to [18].

B. General requests and Poisson distributed servers (GRPS)

From our discussion in Section VI-A1, it is clear that when servers are distributed according to a Poisson process, the exceptional service time distribution equals the regular batch service time distribution. In such a case we have the following queueing model.

Under GRPS, inter-arrival times and batch service times are, respectively, arbitrarily and exponentially distributed. Before initiating a service, a server finds the system in any of the following conditions. (i) $1 \leq n \leq c-1$ and (ii) $n \geq c$. Here n is the number of customers in the waiting buffer. For case (i) the server provides service to all n customers and admits subsequent arrivals until c is reached. For case (ii) the server takes c customers with no admission for subsequent customers arriving within its service time.

In such a case ESABQ can directly be modeled as a special case of a renewal input bulk service queue with accessible and non-accessible batches proposed in [8] with parameter values $a = 1$ and $d = b = c$. Let N_s and N_q denote random variables for numbers of customers in the system and in the waiting buffer respectively for ESABQ under GRPS. We borrow the following definitions from [8].

$$\begin{aligned} P_{n,0} &= \Pr[N_s = n]; 0 \leq n \leq c-1 \\ P_{n,1} &= \Pr[N_q = n]; n \geq 0. \end{aligned} \quad (14)$$

Using results from [8] we obtain the following expressions for equilibrium queue length probabilities.

$$\begin{aligned} P_{0,1} &= \frac{C}{\mu} \left[\frac{r_0^{c-1} - r_0^c}{1 - r_0^c} + \frac{1}{r_0} - 1 \right], \\ P_{n,1} &= \frac{C r_0^{n-1} (1 - r_0)}{\mu (1 - r_0^c)}; n \geq 1, \end{aligned} \quad (15)$$

where $0 < r_0 < 1$ is the real root of the equation $r = F_Y^*(\mu - \mu r^c)$ and C is the normalization constant⁷ given by

$$C = \lambda \left[\frac{1 - \omega^c}{1 - \omega} + \frac{1}{1 - r_0} - \frac{\omega(r_0 - F_Y^*(\mu))}{r_0^c(1 - r_0\omega)} \left(\frac{1 - r_0^c}{1 - r_0} - r_0^{c-1} \frac{1 - \omega^c}{1 - \omega} \right) \right]^{-1}, \quad (16)$$

with $\omega = 1/F_Y^*(\mu)$. We then derive the expected queue length as

$$\begin{aligned} \mathbb{E}[N_q] &= \sum_{n=0}^{\infty} n P_{n,1} = \sum_{n=1}^{\infty} n \frac{C r_0^{n-1} (1 - r_0)}{\mu (1 - r_0^c)} \\ &= \frac{C(1 - r_0)}{\mu (1 - r_0^c)} \sum_{n=1}^{\infty} n r_0^{n-1} = \frac{C}{\mu (1 - r_0^c) (1 - r_0)}. \end{aligned} \quad (17)$$

Applying Little's law and considering the analogy between our service network and ESABQ we obtain the following expression for the expected request distance.

$$\mathbb{E}[D] = \frac{C}{\lambda \mu (1 - r_0^c) (1 - r_0)} + \frac{1}{\mu}. \quad (18)$$

C. Poisson distributed requests and general distributed servers (PRGS)

As discussed in Section VI-A1, if servers are placed on a 1-d line according to a renewal process with requests being Poisson distributed, the service time distribution for the first batch in a busy period differs from those of subsequent batches. Below we derive expressions for queue length distribution and expected request distance for ESABQ under PRGS.

1) *Queue length distribution:* We use a supplementary variable technique to derive the queue length distribution for ESABQ under PRGS as follows.

Let $L(t)$ be the number of customers at time $t \geq 0$, $R(t)$ the residual service time at time $t \geq 0$ (with $R(t) = 0$ if $L(t) = 0$), and $I(t)$ the type of service at time $t \geq 0$ with $I(t) = 1$ (resp. $I(t) = 2$) if exceptional (resp. ordinary) service time.

Let us write the Chapman-Kolmogorov equations for the Markov chain $\{(L(t), R(t), I(t)), t \geq 0\}$.

For $t \geq 0$, $n \geq 1$, $x > 0$, $i = 1, 2$ define

$$\begin{aligned} p_t(n, x; i) &= \mathbb{P}(L(t) = n, R(t) < x, I(t) = i) \\ p_t(0) &= \mathbb{P}(L(t) = 0). \end{aligned}$$

Also, define for $x > 0$, $i = 1, 2$,

$$p(n, x; i) = \lim_{t \rightarrow \infty} p_t(n, x; i) \quad \text{and} \quad p(0) = \lim_{t \rightarrow \infty} p_t(0).$$

By analogy with the analysis for the M/G/1 queue we get

$$\frac{\partial}{\partial t} p_t(0) = -\lambda p_t(0) + \sum_{k=1}^c \frac{\partial}{\partial x} p_t(k, 0; 1) + \sum_{k=1}^c \frac{\partial}{\partial x} p_t(k, 0; 2),$$

so that, by letting $t \rightarrow \infty$,

$$\lambda p(0) = \sum_{k=1}^c \left(\frac{\partial}{\partial x} p(k, 0; 1) + \frac{\partial}{\partial x} p(k, 0; 2) \right). \quad (19)$$

⁷The normalization constant C derived in [8] is incorrect. The correct constant for our case is given in (16).

With further simplification (See [18]), for $n \geq 1, x > 0$ we get

$$\begin{aligned} & \frac{\partial}{\partial x} g(n, x) - \lambda g(n, x) - \frac{\partial}{\partial x} g(n, 0) + \lambda g(n-1, x) \mathbf{1}(n \geq 2) \\ & + \lambda p(0) F_Z(x) \mathbf{1}(n=1) + F_X(x) \frac{\partial}{\partial x} g(n+c, 0) = 0, \end{aligned} \quad (20)$$

where $g(n, x) = p(n, x; 1) + p(n, x; 2)$ for $n \geq 1, x > 0$. Introduce

$$G(z, s) := \sum_{n \geq 1} z^n \int_0^\infty e^{-sx} g(n, x) dx \quad \forall |z| \leq 1, s \geq 0.$$

Denote by $F_Z^*(s) = \int_0^\infty e^{-sx} dF_Z(x)$ the LST of F_Z for $s \geq 0$. Note that

$$\int_0^\infty e^{-sx} F_{Z/X}(x) dx = \frac{F_{Z/X}^*(s)}{s}, \quad \forall s > 0.$$

Multiplying both sides of (20) by $z^n e^{-sx}$, integrating over $x \in [0, \infty)$ and summing over all $n \geq 1$, yields

$$\begin{aligned} (\lambda(1-z) - s) G(z, s) &= \lambda z p(0) F_Z^*(s) - \sum_{n \geq 1} z^n \frac{\partial}{\partial x} g(n, 0) \\ &+ F_X^*(s) \sum_{n \geq 1} z^n \frac{\partial}{\partial x} g(n+c, 0) \end{aligned} \quad (21)$$

where $\lambda p(0) = \sum_{k=1}^c \frac{\partial}{\partial x} g(k, 0)$ from (19). We have

$$\frac{1}{z^c} \sum_{n \geq 1} z^{n+c} \frac{\partial}{\partial x} g(n+c, 0) = \frac{1}{z^c} \sum_{n \geq 1} z^n \frac{\partial}{\partial x} g(n, 0) - \frac{1}{z^c} H(z) \quad (22)$$

where $H(z) = \sum_{k=1}^c z^k a_k$ with $a_k := \frac{\partial}{\partial x} g(k, 0)$, for $k = 1, \dots, c$. Introducing the above into (21) gives

$$\begin{aligned} (\lambda(1-z) - s) G(z, s) &= \left(\frac{F_X^*(s)}{z^c} - 1 \right) \Psi(z) \\ &- F_X^*(s) \frac{H(z)}{z^c} + \lambda z p(0) F_Z^*(s) \end{aligned} \quad (23)$$

where $\Psi(z) := \sum_{n \geq 1} z^n \frac{\partial}{\partial x} g(n, 0)$. Since $G(z, s)$ is well-defined for $|z| \leq 1$ and $s \geq 0$, the r.h.s. of (23) must vanish when $s = \lambda(1-z)$. This gives the relation

$$\Psi(z) = \frac{z^c}{z^c - F_X^*(\theta(z))} \left[-F_X^*(\theta(z)) \frac{H(z)}{z^c} + \lambda z p(0) F_Z^*(\theta(z)) \right]$$

with $\theta(z) = \lambda(1-z)$ and $|z| \leq 1$. Introducing the above in (23) gives

$$\begin{aligned} (\lambda(1-z) - s) G(z, s) &= -F_X^*(s) \frac{H(z)}{z^c} + \lambda z p(0) F_Z^*(s) \\ &+ \frac{F_X^*(s) - z^c}{z^c - F_X^*(\theta(z))} \left[\lambda z p(0) F_Z^*(\theta(z)) - F_X^*(\theta(z)) \frac{H(z)}{z^c} \right]. \end{aligned} \quad (24)$$

Let $N(z)$ be the z -transform of the stationary number of customers in the system. Integrating by part, we get for $n \geq 1$,

$$s \int_0^\infty e^{-sx} g(n, x) dx = \int_0^\infty e^{-sx} dg(n, x),$$

so that

$$\begin{aligned} \lim_{s \rightarrow \infty} s \int_0^\infty e^{-sx} g(n, x) dx &= \lim_{s \rightarrow 0} \int_0^\infty e^{-sx} dg(n, x) \\ &= \int_0^\infty dg(n, x) = g(n, \infty), \end{aligned} \quad (25)$$

where the interchange between the limit and the integral sign is justified by the bounded convergence theorem. Therefore,

$$\begin{aligned} N(z) &= \sum_{n \geq 1} z^n g(n, \infty) + p(0) \\ &= \sum_{n \geq 1} z^n \lim_{s \rightarrow \infty} s \int_0^\infty e^{-sx} g(n, x) dx \quad \text{from (25)} \\ &= \lim_{s \rightarrow 0} s G(z, s) + p(0), \end{aligned} \quad (26)$$

where the interchange between the summation over n and the integral sign is again justified by the bounded convergence theorem. Letting now $s \rightarrow 0$ in (24) and using (26), gives

$$\begin{aligned} \theta(z) N(z) &= \frac{1-z^c}{z^c - F_X^*(\theta(z))} \left[-F_X^*(\theta(z)) \frac{H(z)}{z^c} + \lambda z p(0) F_Z^*(\theta(z)) \right] \\ &- \frac{H(z)}{z^c} + \lambda p(0). \end{aligned} \quad (27)$$

By noting that $\lambda p(0) = \sum_{k=1}^c a_k$ (cf. (19)), Eq. (27) can be rewritten as

$$\begin{aligned} N(z) &= \frac{1}{\theta(z)} \left(\frac{z(1-z^c)}{z^c - F_X^*(\theta(z))} \sum_{k=1}^c a_k \left[F_Z^*(\theta(z)) - z^{k-c-1} F_X^*(\theta(z)) \right] \right. \\ &\left. + \sum_{k=1}^c a_k (1 - z^{k-c}) \right). \end{aligned} \quad (28)$$

The r.h.s. of (28) contains c unknown constants a_1, \dots, a_c yet to be determined. Define $A(z) = F_X^*(\theta(z))$. It can be shown that $z^c - A(z)$ has $c-1$ zeros inside and one on the unit circle, $|z| = 1$ (See [18]). Denote by ξ_1, \dots, ξ_q the $1 \leq q \leq c$ distinct zeros of $z^c - A(z)$ in $\{|z| \leq 1\}$, with multiplicity n_1, \dots, n_q , respectively, with $n_1 + \dots + n_q = c$. Hence,

$$z^c - F_X^*(\theta(z)) = \gamma \prod_{i=1}^q (z - \xi_i)^{n_i}.$$

Since $z^c - A(z)$ vanishes when $z = 1$ and that $\frac{d}{dz}(z^c - A(z))|_{z=1} = c - \rho > 0$, we conclude that $z^c - A(z)$ has one zero of multiplicity one at $z = 1$.

Without loss of generality assume that $\xi_q = 1$ and let us now focus on the zeros ξ_1, \dots, ξ_{q-1} . When $z = \xi_i, i = 1, \dots, q-1$, the term $F_Z^*(\theta(z)) - z^{k-c-1} F_X^*(\theta(z))$ in (28) must have a zero of multiplicity (at least) n_i since $N(\xi_i)$ is well defined. This gives $c-1$ linear equations to be satisfied by ξ_1, \dots, ξ_q . In the particular case where all zeros have multiplicity one (see [18]), namely $q = c$, these $c-1$ equations are

$$\sum_{k=1}^c a_k [F_Z^*(\theta(\xi_i)) - \xi_i^{k-c-1} F_X^*(\theta(\xi_i))] = 0, \quad i = 1, \dots, c-1. \quad (29)$$

With $U(z) := F_Z^*(\theta(z))/F_X^*(\theta(z))$ (29) is equivalent to

$$\sum_{k=1}^c a_k [U(\xi_i) - \xi_i^{k-c-1}] = 0, \quad i = 1, \dots, c-1, \quad (30)$$

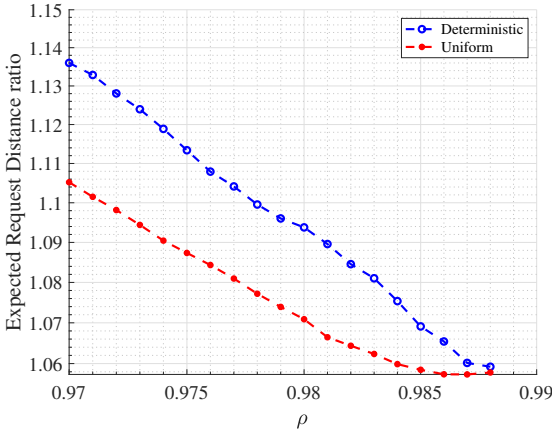


Fig. 3: The plot shows the ratio $\mathbb{E}[D]/\bar{D}_s$ for deterministic and uniform inter-server distance distributions.

since $F_X^*(\theta(\xi_i)) \neq 0$ for $i = 1, \dots, c-1$ ($F_X^*(\theta(\xi_i)) = 0$ implies that $\xi_i = 0$ which contradicts that ξ_i a zero of $z^c - F_X^*(\theta(z))$ since $F_X^*(\theta(0)) = F_X^*(\lambda) > 0$). Eq. (28) can be rewritten as

$$N(z) = \frac{\sum_{k=1}^c a_k [z^c - z^k + z(1 - z^c)F_Z^*(\theta(z)) - (1 - z^k)F_X^*(\theta(z))]}{\theta(z)(z^c - F_X^*(\theta(z)))}. \quad (31)$$

A c -th equation is provided by the normalizing condition $N(z) = 1$. Since the numerator and denominator in (31) have a zero of order 2 at $z = 1$, differentiating twice the numerator and the denominator w.r.t z and letting $z = 1$ gives

$$\sum_{k=1}^c a_k (c(1 + \rho_z) - \rho k) = \lambda(c - \rho), \quad (32)$$

where $\rho_z = \lambda\alpha_Z$. We consider few special cases of the model in [18] and verify with the expressions of queue length distribution available in the literature.

2) *Expected request distance:* From (31) the expected queue length is

$$\begin{aligned} \bar{N} &= \left. \frac{d}{dz} N(z) \right|_{z=1} \\ &= \frac{1}{2\lambda(c - \rho)^2} \sum_{k=1}^c a_k \left[\lambda^2 \sigma_Z^{(2)} c(c - \rho) + \lambda^2 \sigma_X^{(2)} c(1 + \rho_z - k) \right. \\ &\quad \left. + (ck(c - k) + k(k - 1)\rho - c(c - 1))\rho + 2c^2 \rho_z - c(c + 1)\rho_z \rho \right], \quad (33) \end{aligned}$$

where $\sigma_Z^{(2)}$ and $\sigma_X^{(2)}$ are the second order moments of distributions F_Z and F_X respectively. Again by applying Little's law and considering the analogy between our service network with ESABQ we get the following expression for the expected request distance.

$$\mathbb{E}[D] = \bar{N}/\lambda. \quad (34)$$

VII. DISCUSSION OF UNIDIRECTIONAL ALLOCATION POLICIES

In this section we describe generalizations of models and results for unidirectional allocation policies. We first consider the case when inter-user and inter-server distances both have general distributions.

A. Heavy traffic limit for general request and server spatial distributions

Consider the case when the inter-user and inter-server distances each are described by general distributions. We assume server capacity, $c = 1$. As $\rho \rightarrow 1$, we conjecture that the behavior of MTR approaches that of the G/G/1 queue. One argument in favor of our conjecture is the following. As $\rho \rightarrow 1$, the busy cycle duration tends to infinity. Consequently, the impact of the exceptional service for the first customer of the busy period on all other customers diminishes to zero as there is an unbounded increasing number of customers served in the busy period.

It is known that in heavy traffic waiting times in a G/G/1 queue are exponentially distributed and the mean sojourn time is given by $\alpha_X + [(\sigma_X^2 + \sigma_Y^2)/2\alpha_Y(1 - \rho)]$ [9]. We expect the expected request distance to exhibit similar behavior. Thus we have the following conjecture.

Conjecture 1. *At heavy traffic i.e. as $\rho \rightarrow 1$, the expected request distance for the G/G/1 spatial system with $c = 1$ is given by*

$$\mathbb{E}[D] = \alpha_X + \frac{\sigma_X^2 + \sigma_Y^2}{2\alpha_Y(1 - \rho)}. \quad (35)$$

Denote by \bar{D}_s the average request distance as obtained from simulation. We plot the ratio $\mathbb{E}[D]/\bar{D}_s$ across various inter-request and inter-server distance distributions in Figure 3. It is evident that as $\rho \rightarrow 1$, the ratio $\mathbb{E}[D]/\bar{D}_s$ converges to 1 across different inter-server distance distributions.

B. Heterogeneous server capacities under PRGS

We now proceed to analyze a setting where server capacity is a random variable. Assume server capacity \mathcal{C} takes values from $\{1, 2, \dots, c\}$ with distribution $\Pr(\mathcal{C} = j) = p_j, \forall j \in \{1, 2, \dots, c\}$, s.t. $\sum_{j=1}^c p_j = 1$ and $p_c > 0$. We also assume the stability condition $\rho < \bar{C}$ where \bar{C} is the average server capacity. Denote H as the random variable associated with number of requests that traverse through a point just after a server location⁸.

1) *Distribution of H :* Let V denote the number of new requests generated during a service period with $k_v = \Pr(V = v), \forall v \geq 0$. According to the law of total probability, it holds that

$$k_v = \int_0^\infty \Pr(V = v | X = \nu) f_X(\nu) d\nu = \frac{1}{v!} \int_0^\infty e^{-\lambda\nu} (\lambda\nu)^v dF_X(\nu).$$

Then the corresponding generating function $K(z)$ is denoted by $K(z) = \sum_{v=0}^\infty k_v z^v = F_X^*(\lambda(1 - z))$. We now consider an embedded Markov chain generated by H . Denote the corresponding transition matrix as M . Then we have

⁸An analysis for the distribution of number of requests that traverse through any random location would involve the notions of exceptional service and accessible batches.

$$M_{m,l} = \begin{cases} \sum_{i=0}^{c-m} k_i P_{i+m}, & 0 \leq m \leq c, l = 0; \\ \sum_{i=0}^c k_{i+l-m} P_i, & 0 \leq m \leq l, l \neq 0; \\ \sum_{i=m-l}^c k_{i+l-m} P_i, & l+1 \leq m \leq c+l, l \neq 0; \\ 0, & o.w., \end{cases} \quad (36)$$

where $P_i = \sum_{j=i}^c p_j$ and $p_0 = 0$. Let $\pi = [\pi_j, j \geq 0]$ and $N(z) = \sum_{j \geq 0} \pi_j z^j$ denote the steady state distribution and its z -transform respectively. π is obtained out by solving

$$\pi_l = \sum_{m=0}^{\infty} \pi_m M_{m,l}, l = 0, 1, \dots \quad (37)$$

Thus we have for $l \in \mathbb{N}$,

$$\begin{aligned} \pi_0 &= \sum_{m=0}^c \pi_m \sum_{i=0}^{c-m} k_i P_{i+m}, \\ \pi_l &= \sum_{m=0}^l \pi_m \sum_{i=0}^c k_{i+l-m} P_i + \sum_{m=l+1}^{c+l} \pi_m \sum_{i=m-l}^c k_{i+l-m} P_i. \end{aligned} \quad (38)$$

Multiplying by z^l and summing over l gives

$$N(z) = E_\pi + v_1(z) + v_2(z) \quad (39)$$

$$E_\pi = \pi_0 \sum_{i=0}^{c-1} k_i P_{i+1} + \sum_{m=1}^{c-1} \pi_m \sum_{i=m}^{c-1} k_{i-m} P_{i+1} \quad (40)$$

$$v_1(z) = \sum_{l=0}^{\infty} z^l \sum_{m=0}^l \pi_m \sum_{i=0}^c k_{i+l-m} P_i \quad (41)$$

$$v_2(z) = \sum_{l=0}^{\infty} z^l \sum_{m=l+1}^{c+l} \pi_m \sum_{i=m-l}^c k_{i+l-m} P_i. \quad (42)$$

The expressions for $v_1(z)$ and $v_2(z)$ can be further simplified (see [18]) to

$$v_1(z) = N(z) \left\{ \sum_{i=0}^c p_i z^{-i} \left[K(z) - \sum_{j=0}^i k_j z^j \right] + \sum_{i=0}^c k_i z^i \right\} \quad (43)$$

$$\begin{aligned} v_2(z) &= \left[\sum_{m=0}^c z^{-m} \sum_{i=m}^c k_{i-m} P_i \left\{ N(z) - \sum_{j=0}^{m-1} \pi_j z^j \right\} \right] \\ &\quad - N(z) \sum_{i=0}^c k_i z^i. \end{aligned} \quad (44)$$

Combining (39), (43) and (44) yields

$$\begin{aligned} N(z) &= E_\pi + N(z) \left\{ K(z) \sum_{i=0}^c p_i z^{-i} \right\} \\ &\quad - \sum_{j=0}^{c-1} \pi_j \sum_{m=1}^{c-j} z^{-m} \sum_{i=m+j}^c k_{i-(m+j)} P_i. \end{aligned} \quad (45)$$

Thus we obtain

$$N(z) = \frac{E_\pi - \sum_{j=0}^{c-1} \pi_j \sum_{m=1}^{c-j} z^{-m} \sum_{i=m+j}^c k_{i-(m+j)} P_i}{1 - K(z) \sum_{i=0}^c p_i z^{-i}}. \quad (46)$$

Multiplying numerator and denominator by z^c yields

$$N(z) = \frac{z^c E_\pi - \sum_{j=0}^{c-1} \pi_j \sum_{m=1}^{c-j} z^{c-m} \sum_{i=m+j}^c k_{i-(m+j)} P_i}{z^c - K(z) \sum_{i=0}^c p_{c-i} z^i}. \quad (47)$$

To determine $N(z)$, we need to obtain the probabilities $\pi_i, 0 \leq i \leq c-1$. It can be shown that the denominator of (47) has $c-1$ zeros inside and one on the unit circle, $|z| = 1$ (See [18]). As $N(z)$ is analytic within and on the unit circle, the numerator must vanish at these zeros, giving rise to c equations in c unknowns.

Let $\xi_q : 1 \leq q \leq c$ be the zeros of $z^c - K(z) \sum_{i=0}^c p_{c-i} z^i$ in $\{|z| \leq 1\}$. W.l.o.g let $\xi_c = 1$. We have the following $c-1$ equations.

$$E_\pi - \sum_{j=0}^{c-1} \pi_j \sum_{m=1}^{c-j} \xi_q^{-m} \sum_{i=m+j}^c k_{i-(m+j)} P_i = 0, \quad i = 1, \dots, c-1, \quad (48)$$

A c -th equation is provided by the normalizing condition $\lim_{z \rightarrow 1} N(z) = 1$. In the particular case where all zeros have multiplicity one, it can be shown that these c equations are linearly independent⁹. Once the parameters $\{\pi_i, 0 \leq i \leq c-1\}$ are known, $\mathbb{E}[H]$ can be expressed as

$$\mathbb{E}[H] = \bar{H} = \lim_{z \rightarrow 1} N'(z). \quad (49)$$

2) *Expected Request Distance*: To evaluate the expected request distance we adopt arguments from [5]. Consider any interval of length ν between two consecutive servers. There are on average \bar{H} requests at the beginning of the interval, each of which must travel ν distance. New users are spread randomly over the interval and there are on an average $\lambda\nu$ new users. The request made by each new user must travel on average $\nu/2$. Thus we have

$$\begin{aligned} \mathbb{E}[D] &= \frac{1}{\rho} \int_0^\infty (\bar{H}\nu + \frac{1}{2}\lambda\nu^2) dF_X(\nu) \\ &= \frac{1}{\rho} \left[\frac{\bar{H}}{\mu} + \frac{\lambda}{2} \left(\sigma_X^2 + \frac{1}{\mu^2} \right) \right]. \end{aligned} \quad (50)$$

C. Uncapacitated request allocation

An interesting special case of the unidirectional general matching is the uncapacitated scenario. Consider the case where servers do not have any capacity constraints, i.e. $c = \infty$. In such a case, all users are assigned to the nearest server to their right.

GRPS: When $c \rightarrow \infty$ and given $0 < r_0 < 1$, $r_0 = F_Y^*(\mu - \mu r_0^c) = F_Y^*(\mu)$. Setting $\omega = 1/F_Y^*(\mu) = 1/r_0$ in (16) and simplifying yields

$$C \rightarrow 0, \text{ as } c \rightarrow \infty, \implies \mathbb{E}[D] \rightarrow \frac{1}{\mu} \text{ as } c \rightarrow \infty. \quad (51)$$

PRGS: Under PRGS, when $c \rightarrow \infty$ there exists no request allocated to a server other than the nearest server to its right.

⁹For all cases evaluated across uniform, deterministic and hyperexponential distributions we found the set of c equations to be linearly independent.

Again using Bailey's method as in [5] and setting $\bar{H} = 0$ in (50) we get

$$\mathbb{E}[D] \rightarrow \frac{\mu}{2} \left(\sigma_X^2 + \frac{1}{\mu^2} \right) \text{ as } c \rightarrow \infty. \quad (52)$$

VIII. BIDIRECTIONAL ALLOCATION POLICIES

Both UGS and MTR minimize expected request distance among all unidirectional policies. In this section we formulate the bi-directional allocation policy that minimizes expected request distance. Let $\eta : R \rightarrow S$ be any mapping of users to servers. Our objective is to find a mapping $\eta^* : R \rightarrow S$, that satisfies

$$\begin{aligned} \eta^* &= \arg \min_{\eta} \sum_{i \in R} d_{\mathcal{L}}(r_i, s_{\eta(i)}) \\ \text{s.t. } \sum_{i \in R} \mathbb{1}_{\eta(i)=j} &\leq c, \forall j \in S \end{aligned} \quad (53)$$

W.l.o.g, let $r_1 \leq r_2 \leq \dots \leq r_i \leq \dots \leq r_{|R|}$ be locations of requests and $s_1 \leq s_2 \leq \dots \leq s_i \leq \dots \leq s_{|S|}$ be locations of servers. We first focus on the case when $c = 1$. We consider the following two scenarios.

Case 1: $|R| = |S|$

When $|R| = |S|$, an optimal allocation strategy is given by the following theorem [6].

Theorem 2. *When $|R| = |S|$, an optimal assignment is obtained by the policy: $\eta^*(i) = i$, $\forall i \in \{1, \dots, |R|\}$ i.e. allocating the i^{th} request to the i^{th} server and the average request distance is given by*

$$\mathbb{E}[D] = \frac{1}{|R|} \sum_{i=1}^{|R|} |s(i) - r(i)|. \quad (54)$$

Case 2: $|R| < |S|$ This is the case where there are fewer requesters than servers. In this case, a Dynamic Programming (DP) based algorithm (Algorithm 1) obtains the optimal assignment.

Let $C[i, j]$ denote the optimal cost (i.e., sum of distances) of assigning the first i requests (counting from the left) located at $r_1 \leq r_2 \leq \dots \leq r_i$ to the first j servers (also counting from the left) located at $s_1 \leq s_2 \leq \dots \leq s_j$. If $j = i$, the optimal assignment is trivial due to Theorem 2 and $C[i, i]$ is computed easily for all $i \leq |R|$ by summing pairwise distances $d[1, 1], d[2, 2], \dots, d[i, i]$ (Lines 6–7). For the base case, $i = 1, j > 1$, only the first user needs to be assigned to its nearest server (Lines 9–16). For the general dynamic programming step, consider $j > i$. Then $C[i, j]$ can be expressed in terms of the costs of two subproblems, i.e., $C[i-1, j-1]$ and $C[i, j-1]$ (Lines 19–24). In the optimal solution, two cases are possible: either request i is assigned to server j , or the latter is left unallocated. The former case occurs if the first $i-1$ requests are assigned to the first $j-1$ servers at cost $C[i-1, j-1]$, and the latter case occurs when the first i requests are assigned to the first $j-1$ servers at cost $C[i, j-1]$. This is a consequence of the no-crossing lemma (Lemma 1). The optimal $C[i, j]$ is

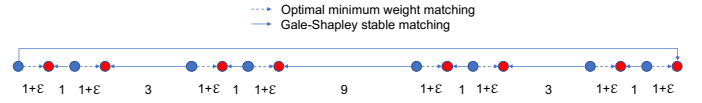


Fig. 4: Worst case scenario for Gale-Shapley.

chosen depending on these two costs and the current distance $d[i, j]$.

Lemma 1. *In an optimal solution, η^* , to the problem of matching users at $r_1 \leq r_2 \leq \dots \leq r_{|R|}$ to servers at $s_1 \leq s_2 \leq \dots \leq s_{|S|}$, where $|S| \geq |R|$, there do not exist indices i, j such that $\eta^*(i) > \eta^*(i')$ when $i' > i$.*

Proof. See [18]. \square

The dynamic programming algorithm fills cells in an $|R| \times |S|$ matrix C whose origin is in the north-west corner. The lower triangular portion of this matrix is invalid since $|R| \leq |S|$. The base cases populate the diagonal and the northernmost row, and in the general DP step, the value of a cell depends on the previously computed values in the cells located to its immediate west and diagonally north-west. As an optimization, for a fixed i , the j -th loop index needs to run only from $i+1$ through $i+|S|-|R|$ (Lines 11 and 18) instead of from $i+1$ through $|S|$. This is because the first request has to be assigned to a server s_j with $j \leq |S| - |R| + 1$ so that the rest of the $|R| - 1$ requests have a chance of being placed on unique servers¹⁰. The optimal average request distance is given by $C[|R|, |S|]$.

The time complexity of the main DP step is $O(|R| \times (|S| - |R| + 1))$. Note that this assumes that the pairwise distance matrix d of dimension $|R| \times |S|$ has been precomputed. The optimization applied above can be similarly applied to this computation and hence the overall time complexity of Algorithm 1 is $O(|R| \times (|S| - |R| + 1))$. Therefore, if $|S| = O(|R|)$, the worst case time complexity is quadratic in $|R|$. However, if $|S| - |R|$ grows only sub-linearly with $|R|$, the time complexity is sub-quadratic in $|R|$.

Note that retrieving the optimal assignment requires more book-keeping. An $|R| \times |S|$ matrix A stores key intermediate steps in the assignment as the DP algorithm progresses (Lines 8, 16, 21, 24). The optimal assignment vector π can be retrieved from matrix A using procedure READOPTASSIGNMENT.

Another bidirectional assignment scheme is the Gale-Shapley algorithm [7], which produces stable assignments, though in the worst case it can yield an assignment that is $O(|R|^{\ln 3/2}) \approx O(|R|^{0.58})$ times costlier than the optimal assignment yielded by Algorithm 1, where $|R|$ is the number of users [19]. The worst case scenario is illustrated in Figure 4, with $|R| = 2^{t-1}$, where t is the number of clusters of users and servers; and the largest distance between adjacent points is

¹⁰Note that in this exposition, we consider server capacity $c = 1$. If $c > 1$, we simply add c servers at each prescribed server location, and requests will still be placed on unique servers.

Algorithm 1 Optimal Assignment by Dynamic Programming

```

1: Input:  $r_1 \leq \dots \leq r_{|R|}$ ;  $s_1 \leq \dots \leq s_{|S|}$ 
2: Output: The optimal assignment  $\pi$ 
3: procedure OPTDP( $r, s$ )
4:    $d_{|R| \times |S|} = \text{COMPUTEPAIRWISEDISTANCES}(r, s)$ 
5:    $C = \{\infty\}_{|R| \times |S|}$ 
6:   for  $i = 1, \dots, |R|$  do
7:      $C[i, i] = \text{TRIVIALASSIGNMENT}(i, d)$ 
8:    $A[|R|, |R|] = |R|$ 
9:    $\text{nearest} = 0$ 
10:   $\text{nearestcost} = C[1, 1]$ 
11:  for  $j = 2, \dots, |S| - |R| + 1$  do
12:    if  $d[1, j] < \text{nearestcost}$  then
13:       $\text{nearestcost} = d[1, j]$ 
14:       $\text{nearest} = j$ 
15:     $C[1, j] = \text{nearestcost}$ 
16:     $A[1, j] = \text{nearest}$ 
17:  for  $i = 2, \dots, |R|$  do
18:    for  $j = i + 1, \dots, |S| - |R| + 1$  do
19:      if  $C[i, j - 1] < d[i, j] + C[i - 1, j - 1]$  then
20:         $C[i, j] = C[i, j - 1]$ 
21:         $A[i, j] = A[i, j - 1]$ 
22:      else
23:         $C[i, j] = d[i, j] + C[i - 1, j - 1]$ 
24:         $A[i, j] = j$ 
25:  return READOPTASSIGNMENT( $A$ )
26: procedure TRIVIALASSIGNMENT( $n, d$ )
27:    $\text{Cost} = 0$ 
28:   for  $i = 1, \dots, n$  do
29:      $\text{Cost} = \text{Cost} + d[i, i]$ 
30:   return  $\text{Cost}$ 
31: procedure READOPTASSIGNMENT( $A$ )
32:    $|R|, |S| = \text{DIMENSIONS}(A)$ 
33:    $s = |S|$ 
34:   for  $i = |R|, \dots, 1$  do
35:      $\pi[i] = A[i, s]$ 
36:      $s = A[i, s] - 1$ 
37:   return  $\pi$ 

```

3^{t-2} . However at low/moderate loads for the cases evaluated in Section IX, we find its performance to be not much worse than optimal.

IX. NUMERICAL EXPERIMENTS

In this section, we examine the effect of various system parameters on expected request distance under MTR policy. We also compare the performance of various greedy allocation strategies along with the unidirectional policies to the optimal strategy.

A. Experimental setup

In our experiments, we consider a mean requester rate $\lambda \in (0, 1)$. We consider various inter-server distance distributions with density one. In particular, (i) for exponential

distributions, the density is set to $\mu = 1$; (ii) for deterministic distributions, we assign parameter $d_0 = 1$. (iii) for second order hyper-exponential distribution (H_2), denote p_1 and p_2 as the phase probabilities. Let μ_1 and μ_2 be corresponding phase rates. We assume $p_1/\mu_1 = p_2/\mu_2$. We express H_2 parameters in terms of the squared coefficient of variation, c_v^2 , and mean inter-server distance, α_X , i.e. we set $p_1 = (1/2)(1 + \sqrt{(c_v^2 - 1)/(c_v^2 + 1)})$, $p_2 = 1 - p_1$, $\mu_1 = 2p_1/\alpha_X$ and $\mu_2 = 2p_2/\alpha_X$. Unless specified, for H_2 we take $c_v^2 = 4$ with $c = 2$. Also if not specified, users are distributed according to a Poisson process and servers according to a renewal process.

We consider a collection of 10^5 users and 10^5 servers, i.e. $|R| = |S| = 10^5$. We assign users to servers according to MTR. Let $R_M \subseteq R$ be the set of users allocated under MTR. Clearly $|R_M| \leq |R|$. We then run optimal and other greedy policies on the set R_M and S . For each of the experiments, the expected request distance for the corresponding policy is averaged over 50 trials.

B. Sensitivity analysis

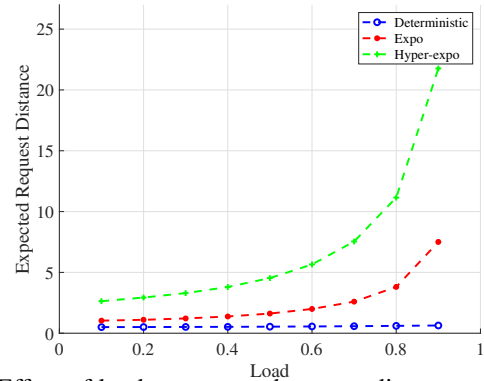


Fig. 5: Effect of load on expected request distance with $c = 2$.

1) *Expected request distance vs. load:* We first study the effect of load ($= \lambda/c\mu$) on $\mathbb{E}[D]$ as shown in Figure 5. Clearly $\mathbb{E}[D]$ increases with load. Note that H_2 distribution exhibits the largest expected request distance and the deterministic distribution, the smallest because the servers are evenly spaced. While for H_2 , c_v^2 is larger than for the exponential distribution. Consequently servers are clustered, which increases $\mathbb{E}[D]$.

2) *Expected request distance vs. squared coefficient of variation:* We now examine how c_v^2 affects $\mathbb{E}[D]$ when ρ is fixed. We compare two systems: a general request with Poisson distributed servers (H_2/M) and a Poisson request with general distributed servers (M/H_2) where the general distribution is a H_2 distribution with the same set of parameters, i.e. we fix $\lambda = \mu = 1$ with $c = 2$. The results are shown in Figure 6. Note that, when $c_v^2 = 1$ H_2 is an exponential distribution and both H_2/M and M/H_2 are identical $M/M/1$ systems. As discussed for the previous graph, performance of both systems decreases with increase in c_v^2 due to increase in the variability of user and server placements. However, from Figure 6 it is clear that performance is more sensitive to server placement as compared to the corresponding user placement.

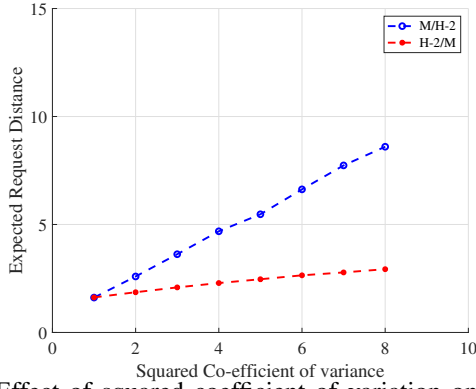


Fig. 6: Effect of squared coefficient of variation on expected request distance with $\lambda = \mu = 1$ and $c = 2$.

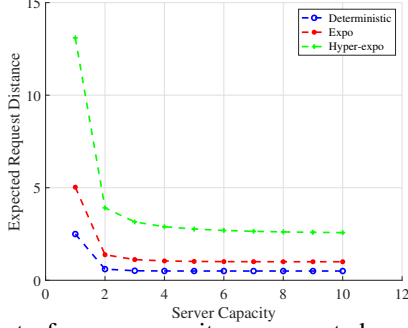


Fig. 7: Effect of server capacity on expected request distance with $\rho = 0.8$.

3) *Expected request distance vs. server capacity:* We now focus on how server capacity affects $\mathbb{E}[D]$ as shown in Figure 7. We fix $\rho = 0.8$. With an increase in c , while keeping ρ fixed, $\mathbb{E}[D]$ decreases. This is because queuing delay decreases. Note that $\mathbb{E}[D]$ gradually converge to a value with increase in server capacity. Theoretically, this can be explained by our discussion on uncapacitated allocation in Section VII-C. As $c \rightarrow \infty$ the contribution of queuing delay to $\mathbb{E}[D]$ vanishes and $\mathbb{E}[D]$ becomes insensitive to c .

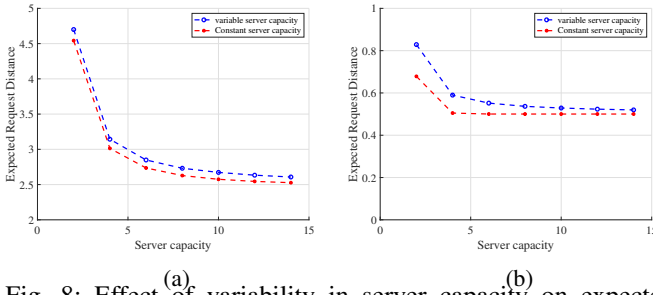


Fig. 8: Effect of variability in server capacity on expected request distance for H_2 (a) and Deterministic (b) distributions with $\rho = 0.8$.

4) *Expected request distance vs. capacity moments:* We investigate the heterogeneous capacity scenario as discussed in Section VII-B. Consider the plot shown in Figure 8. We fix $\rho = 0.8$. For the variable server capacity curve we choose a value for server capacity for each server uniformly at random from the set $\{1, 2, \dots, 2c\}$. For the constant server capacity curve we deterministically assign server capacity c to each

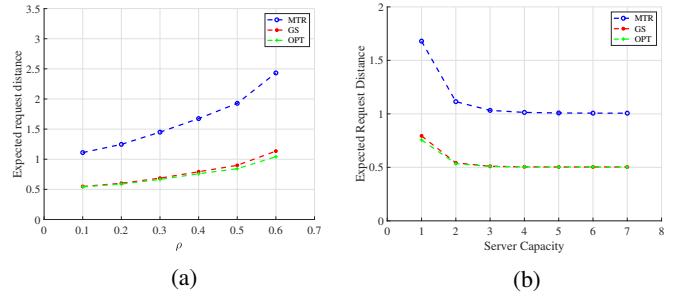


Fig. 9: Comparison of different allocation policies: (a) ρ vs $\mathbb{E}[D]$ with $c = 1$, (b) c vs. $\mathbb{E}[D]$ with $\rho = 0.4$.

server. While both the curves exhibit similar performance under H_2 distribution, we observe better performance for constant server capacity curve at lower values of c under Deterministic distribution. Variability in constant server case is zero, thus explaining its better performance.

C. Comparison of different allocation policies

We consider the case in which both users and servers are distributed according to Poisson processes. From Figure 9 (a), we observe that due to its directional nature MTR has a larger expected request distance compared to other policies while GS provides near optimal performance. In Figure 9 (b), we compare the performance of allocation policies across different server capacities. The expected request distance decreases with increase in server capacities across all policies. Both GS and the optimal policy converge to the same value as c gets higher.

We observe similar trends in the case of deterministic inter-server distance distributions. However, under equal densities, all the policies produce smaller expected request distance as compared to their Poisson counterpart. This advocates for placing equidistant servers in a bidirectional system with Poisson distributed requesters to minimize expected request distance.

X. CONCLUSION

We introduced a queuing theoretic model for analyzing the behavior of unidirectional policies to allocate tasks to servers on the real line. We showed the equivalence of UGS and MTR w.r.t the expected request distance and presented results associated with the case when either requesters or servers are Poisson distributed. In this context, we analyzed a new queuing theoretic model: ESABQ, not previously studied in queuing literature. We also proposed a dynamic programming based algorithm to obtain an optimal allocation policy in a bi-directional system. We performed sensitivity analysis for unidirectional system and compared the performance of various greedy allocation strategies along with the unidirectional policies to that of optimal policy. Going further, we aim to extend our analysis for unidirectional policies to a two-dimensional geographic region.

XI. ACKNOWLEDGMENT

This research was sponsored by the U.S. ARL and the U.K. MoD under Agreement Number W911NF-16-3-0001

and by the NSF under Grant CNS-1617437. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the National Science Foundation, U.S. ARL or the U.K. MoD. This document does not contain technology or technical data controlled under either the U.S. International Traffic in Arms Regulations or the U.S. Export Administration Regulations. The authors thank Prof. Amotz Bar-Noy for initial discussions on the optimal dynamic programming solution.

REFERENCES

- [1] H. K. Abadi and B. Prabhakar. Stable Matchings in Metric Spaces: Modeling Real-World Preferences using Proximity. *arXiv:1710.05262*, 2017.
- [2] P. Agarwal, A. Efrat, and M. Sharir. Vertical Decomposition of Shallow Levels in 3-Dimensional Arrangements and Its Applications. *SOCG*, 1995.
- [3] R. Ahuja, T. Magnanti, and J. Orlin. *Network Flows: Theory, Algorithms, and Applications*. Prentice-Hall, Inc, 1993.
- [4] L. Atzori, A. Iera, and G. Morabito. The Internet of Things: A Survey. *Computer Networks*, 54(15):2787–2805, 2010.
- [5] N. T. J. Bailey. On Queueing Processes with Bulk Service. *J. R. Stat. SOCE.*, 16:80–87, 1954.
- [6] J. Bukac. Matching On a Line. *arXiv:1805.00214*, 2018.
- [7] D. Gale and L. Shapley. College Admissions and Stability of Marriage. *Amer. Math. Monthly* 69, pages 9–15, 1962.
- [8] V. Goswami and P. V. Laxmi. A Renewal Input Single and Batch Service Queues with Accessibility to Batches. *International Journal of Management Science and Engineering Management*, pages 366–373, 2011.
- [9] D. Gross and C. Harris. Fundamentals of Queueing Theory. *Wiley Series in Probability and Statistics*, 1998.
- [10] I. W. H. Ho, K. K. Leung, and J. W. Polak. Stochastic Model and Connectivity Dynamics for VANETs in Signalized Road Systems. *IEEE/ACM Transactions on Networking*, 19(1):195–208, 2011.
- [11] A. E. Holroyd, R. Pemantle, R. Peres, and O. Schramm. Poisson Matching. *Annales de l IHP Probabilites et Statistiques*, 45:266–287, 2009.
- [12] I. Jawhar, N. Mohamed, and P. D. Agrawal. Linear Wireless Sensor Networks: Classification and Applications. *Journal of Network and Computer Applications*, 34(5):1671–1682, 2011.
- [13] F. Kingman. The Effect of Queue Discipline on Waiting Time Variance. *Math. Proc. Cambridge Phil. Soc.*, 58:163–164, 1962.
- [14] L. Kleinrock. *Queueing Systems*. John Wiley and Sons, 1976.
- [15] K. K. Leung, W. A. Massey, and W. Whitt. Traffic Models for Wireless Communication Networks. *IEEE Journal on Selected Areas in Communications*, 12(8):1353–1364, 1994.
- [16] M. Mezard and G. Parisi. The Euclidean Matching Problem. *J. Phys. France*, 49:2019–2025, 1988.
- [17] J. Orlin. A Polynomial Time Primal Network Simplex Algorithm for Minimum Cost Flows. *Mathematical Programming*, 78:109–129, 1997.
- [18] N. K. Panigrahy, P. Basu, P. Nain, , D. Towsley, A. Swami, K. S. Chan, and K. K. Leung. Resource Allocation in One-dimensional Distributed Service Networks. *Arxiv preprint arXiv:1901.02414*, 2019.
- [19] E. M. Reingold and R. E. Tarjan. On a Greedy Heuristic for Complete Matching. *SIAM Journal on Computing*, 10(4):676–681, 1981.
- [20] P. Welch. On a Generalized m/g/1 Queueing Process in Which The First Customer of Each Busy Period Receives Exceptional Service. *Operations Research*, 12:736–752, 1964.