This article was downloaded by: [143.215.16.100] On: 28 December 2023, At: 13:57 Publisher: Institute for Operations Research and the Management Sciences (INFORMS) INFORMS is located in Maryland, USA



Operations Research

Publication details, including instructions for authors and subscription information: http://pubsonline.informs.org

Dynamic Pricing and Matching for Two-Sided Queues

Sushil Mahavir Varma, Pornpawee Bumpensanti, Siva Theja Maguluri, He Wang

To cite this article:

Sushil Mahavir Varma, Pornpawee Bumpensanti, Siva Theja Maguluri, He Wang (2023) Dynamic Pricing and Matching for Two-Sided Queues. Operations Research 71(1):83-100. https://doi.org/10.1287/opre.2021.2233

Full terms and conditions of use: https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2022, INFORMS

Please scroll down for article-it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit http://www.informs.org

Vol. 71, No. 1, January-February 2023, pp. 83-100 ISSN 0030-364X (print), ISSN 1526-5463 (online)

https://pubsonline.informs.org/journal/opre

Methods

Dynamic Pricing and Matching for Two-Sided Queues

Sushil Mahavir Varma,^{a,}* Pornpawee Bumpensanti,^a Siva Theja Maguluri,^a He Wang^a

^aSchool of Industrial and Systems Engineering, Georgia Institute of Technology, Atlanta, Georgia 30332 *Corresponding author

Contact: sushil@gatech.edu, 📵 https://orcid.org/0000-0001-7855-3447 (SMV); pornpawee@gatech.edu (PB); siva.theja@gatech.edu (STM); he.wang@isye.gatech.edu, https://orcid.org/0000-0001-7444-2053 (HW)

Received: February 25, 2020 Revised: February 12, 2021; July 10, 2021

Accepted: October 13, 2021 Published Online in Articles in Advance: February 8, 2022

Area of Review: Stochastic Models

https://doi.org/10.1287/opre.2021.2233

Copyright: © 2022 INFORMS

Abstract. Motivated by applications from gig economy and online marketplaces, we study a two-sided queueing system under joint pricing and matching controls. The queueing system is modeled by a bipartite graph, where the vertices represent customer or server types and the edges represent compatible customer-server pairs. Both customers and servers sequentially arrive to the system and join separate queues according to their types. The arrival rates of different types depend on the prices set by the system operator and the expected waiting time. At any point in time, the system operator can choose certain customers to match with compatible servers. The objective is to maximize the long-run average profit for the system. We first propose a fluid approximation-based pricing and maximum-weight (max-weight) matching policy, which achieves an $O(\sqrt{\eta})$ optimality rate when all the arrival rates are scaled by η . We further show that a two-price and maxweight matching policy achieves an improved $O(\eta^{1/3})$ optimality rate. Under a broad class of pricing policies, we prove that any matching policy has an optimality rate that is lower bounded by $\Omega(\eta^{1/3})$. Thus, the latter policy achieves the optimal rate with respect to η . We also demonstrate the advantage of max-weight matching with respect to the number of server and customer types n. Under a complete resource pooling condition, we show that max-weight matching achieves $O(\sqrt{n})$ and $O(n^{1/3})$ optimality rates for static and two-price policies, respectively, and the latter matches the lower bound $\Omega(n^{1/3})$. In comparison, the randomized matching policy may have an $\Omega(n)$ optimality rate.

Funding: This work was supported in part by the National Science Foundation [Grant CMMI-2145661]. Supplemental Material: The online appendix is available at https://doi.org/10.1287/opre.2021.2233.

Keywords: queueing • dynamic pricing • max-weight matching • Markov decision process

Downloaded from informs org by [143.215.16.100] on 28 December 2023, at 13:57. For personal use only, all rights reserved

Most queueing models consider a fixed set of servers with sequentially arriving customers. In this paper, however, we consider a two-sided queueing system where servers also arrive sequentially and then wait to be matched with customers. Various applications of online marketplaces and gig economy platforms can be modeled as two-sided queues—for example, Uber and Lyft, where passengers are matched with drivers; Uber Eats and DoorDash, where customer orders are matched with meal delivery couriers; and crowdsourced workforce platforms, such as TaskRabbit, where tasks are matched with contributors. Most of these platforms use both dynamic pricing and dynamic matching as levers to control the marketplaces.

Motivated by these applications, we consider a canonical model of two-sided queues with multiple types of servers and customers. Each customer type is compatible with a subset of server types. For example, in the case of ride-hailing marketplaces, the types of servers (drivers) and customers are determined by the

proximity of their current locations as well as other factors such as the numbers of seats requested by passengers and the vehicle capacities. Our model assumes a fairly general setting with arbitrary numbers of customer and server types, with their compatibility modeled by a bipartite graph.

At each point in time, the system operator posts a price for each customer and server type. Then, customers and servers who are willing to accept the quoted prices (after they factor in expected waiting costs) will enter the system. Those who entered will wait in queues separated by their types until they are matched to a compatible counterpart type. After a customerserver pair is matched, the pair will leave the queueing system immediately to complete the service. The system operator earns a profit that is equal to the difference between the price charged to the customer and the price quoted to the server.

We formulate the system as a Markov decision process (MDP) in the infinite time horizon. The operator can vary the prices for different customer and server types as well as decide when to match and which customer-server pair to match. The objective is to maximize the long-run average profit obtained by the system operator.

There are several technical challenges to analyzing such a stochastic system. The first challenge is the curse of dimensionality in solving and analyzing the MDP. As the number of customer or server types increases, the dimension of the state space increases exponentially (even when the buffer size of each queue is bounded). It is hence intractable to solve the exact MDP for large-scale systems with multiple types. In this paper, we propose several approximate policies to obtain near-optimal solutions for the MDP.

The second challenge is that the stochastic behavior of the two-sided queueing system is complicated by the interplay between pricing and matching decisions. Our proposed policies use dynamic pricing to ensure the stability of the two-sided queue system, so the arrival rates of customers and servers vary with the queue lengths. As the queue lengths change, the matching decisions among different types are adjusted dynamically, which in turn, affects the system state and pricing decisions. As a result, the queue lengths of different types are intricately correlated. The system cannot be decomposed into a set of simple queue, and the pricing and matching decisions cannot be decoupled and analyzed separately. To solve this challenge, we use the Lyapunov drift method to analyze the stochastic system as a whole in order to bound the total queue length.

1.1. Summary of Results

We first present a fluid model for the two-sided queueing system and show that the profit obtained by the fluid model is an upper bound on the achievable profit under any policy. Based on the fluid model, we propose several pricing and matching policies.

In Section 3, we analyze the proposed policies in a large-scale regime in which the arrival rates of all types are scaled by a factor $\eta \to \infty$. We consider a static pricing policy using the fluid solution combined with the maximum-weight (max-weight) matching algorithm. We show that the profit loss of this policy from the fluid solution benchmark is $O(\sqrt{\eta})$ (Section 3.2). We then propose a generalization of the fluid pricing policy that uses two prices for each queue type (see Kim and Randhawa 2017). For the two-price policy combined with max-weight matching policy, we show that the profit loss from the fluid solution benchmark is reduced to $O(\eta^{1/3})$ (Section 3.3). Furthermore, we prove that for a broad class of pricing policies, using any matching policy will result in a profit loss lower bounded by $\Omega(\eta^{1/3})$ (Section 3.4).

In Section 4, we consider a *large-system* regime in which both the number of server/customer types and

the arrival rates are scaled $(n \to \infty, \eta \to \infty)$. We show that the max-weight algorithm is *delay optimal*. In particular, max-weight matching minimizes the revenue loss under fluid pricing and two-price policies among all matching policies (Section 4.1). Under the complete resource pooling condition, we characterize the profit loss of max-weight matching; the profit loss scales as $O(\sqrt{n\eta})$ for the fluid pricing policy and $O((n\eta)^{1/3})$ for the two-price policy (Section 4.2). Furthermore, we establish a lower bound showing that any pricing and matching will incur a profit loss of $\Omega((n\eta)^{1/3})$, so the two-price max-weight policy is asymptotically optimal (Section 4.3). In contrast, if one directly applies the solution of the fluid model as a state-independent randomized matching policy, the profit loss scales as $O(n\eta^{1/2})$ for the fluid pricing policy and $O(n\eta^{1/3})$ for the two-price policy.

In Online Appendix A, we further analyze the structure of the MDP model and propose approximate Dynamic Programming solutions. In some special cases, we are able to show structural properties of the optimal dynamic pricing policy. In addition, we present an Linear Programming-based approximation technique with a constraint generation algorithm to solve the MDP efficiently.

1.2 Literature Review

1.2.1. Dynamic Matching. Dynamic matching markets have numerous applications, such as ride sharing (Banerjee et al. 2017), e-commerce marketplaces like Amazon.com or eBay, kidney exchange (Roth et al. 2007, Anderson et al. 2017), and payment processing networks (Sivaraman et al. 2020). We will discuss previous work involving dynamic matching in the context of two-sided queues.

Caldentey et al. (2009) and Adan and Weiss (2012) considered bipartite matching for two-sided queues on a first come, first served (FCFS) basis; each arriving customer is matched to a compatible server who has the earliest arrival time and has not been matched. Under this matching rule, they analyzed steady-state matching rates between certain customer and server types. Furthermore, they deduced the necessary conditions on the frequency of arrivals for stability of the system and also derived the stationary distribution. Gurvich and Ward (2014) analyzed a general multisided queuing system, which includes the two-sided queueing system as a special case. Their objective was to minimize holding cost in a finite horizon. They presented a periodic review matching algorithm and showed asymptotic optimality as arrival rates become large.

Hu and Zhou (2021) studied a two-sided matching system similar to ours. Their goal is to maximize the discounted reward obtained by matching customers and servers in a finite horizon while accounting for the holding costs. They study conditions such that a priority rule is optimal. In addition, they present a matching algorithm based on fluid approximation and show that it is asymptotically optimal. The main distinction of Hu and Zhou (2021) with our paper is that they do not consider dynamic pricing. In addition, although they use fluid approximation to generate static (open-loop) matching decisions, we use the maxweight algorithm to generate (closed-loop) matching decisions that are adaptive to queue lengths. Chen and Hu (2020) studied a dynamic pricing and matching problem in which strategic customers and servers arrive dynamically and have heterogeneous waiting costs. Their paper assumes that all customers and servers are compatible and considers a greedy matching policy on a first come, first served basis.

Dynamic matching problems were also studied in the context of kidney exchanges, albeit in a nontwo-sided setting in Anderson et al. (2017) and Akbarpour et al. (2020). Pricing is usually forbidden in kidney exchanges because of ethical and legal reasons. These papers study the value of "batching" (i.e., holding compatible matching pairs in the hope that better matching will arrive in future). However, both papers find that batching in general does not provide significant benefit.

1.2.2. Dynamic Pricing for Queues. First, we discuss the literature involving dynamic pricing in the context of single-sided queues and then review those involving two-sided queues.

Low (1974a) is one of the earlier works studying dynamic pricing in a single-sided queue. The paper considered price-dependent customer arrivals with a finite buffer; the rewards include the payment by customers and holding costs incurred by the operator. Monotonicity of the optimal pricing policy is showed. It was later extended to infinite buffer capacity in Low (1974b). Chen and Frank (2001) considered a queuing model with customers who are sensitive to both waiting time and price. They presented structural properties on optimal pricing decisions and monotonicity of optimal bias function. In the context of network services like call centers, Paschalidis and Tsitsiklis (2000) considered a system with finite total resource. They consider price dependent customer arrivals belonging to diverse types differing in resource requirements. The objective is to find a pricing policy to maximize revenue. They show multiple structural properties, like concavity of value function and monotonicity of optimal policy.

Kim and Randhawa (2017) considers a single-server queuing system and studies the benefit of dynamic pricing over static pricing. They assume that the customers are delay sensitive and consider a revenue maximization objective. They present a static pricing policy and a two-price policy, and they also provide the rate of convergence of these policies. Our two-

price policy considered in Section 3.3 is motivated by the results from Kim and Randhawa (2017). The method of Kim and Randhawa (2017) involves applying the Taylor series expansion to the revenue function and then, bounding the expected steady-state queue length. The main distinction of Kim and Randhawa (2017) with our paper is that they consider a singleserver queue, whereas we consider a network of twosided queues with matching decisions. It is nontrivial to generalize the method presented in Kim and Randhawa (2017) to a two-sided queueing network, as an exact analysis of the steady-state distribution is intractable because of the complex interaction among different queues. In addition, unlike the single-server setting in Kim and Randhawa (2017), matching decisions play a critical role in our model and cannot be decoupled from the pricing decisions. Aside from establishing asymptotic rates with large arrival rates, we also complement the result in Kim and Randhawa (2017) by showing the advantage of twoprice pricing (when combined with appropriate matching policies) for large network sizes.

The joint problem of dynamic pricing and matching was also studied by Özkan and Ward (2020) under the objective of maximizing the number of successful matches. They proposed an asymptotically optimal pricing and matching policy with large arrival rates. The differences with our work are that they proposed static policies based on the fluid model and analyzed the system for a finite time horizon.

A two-sided queueing model with both customer and server arrivals is studied by Nguyen and Stolyar (2018). They consider a setting where the arrival rate of the servers can be controlled. However, the focus in Nguyen and Stolyar (2018) was to establish system stability and process-level convergence, whereas the objective in our model is to maximize profit.

Several recent papers have studied dynamic pricing in the context of ride-hailing systems (Besbes et al. 2021, Yan et al. 2020, Hu et al. 2021). Banerjee et al. (2017, 2018) studied a closed queuing network, where the number of cars in the system is a constant and the customers abandon the system if they are not matched immediately. Banerjee et al. (2017) considered a state-independent pricing policy and prove the approximation ratio with respect to optimal pricing policy. Banerjee et al. (2016) proposed a state-dependent pricing policy and argue that the benefit of dynamic pricing is in the robustness of the performance of the system.

In sum, most of the previous work on dynamic matching either is in the context of single-sided queues or is not coupled with revenue optimization. Of the few that consider both of these, the matching policy considered is an open-loop policy. On the other hand, we consider all of these aspects and show the asymptotic optimality under closed-loop matching policies.

1.2.3. Max-Weight Algorithm. In this work, we apply a max-weight matching algorithm to two-sided queuing systems. This algorithm was first proposed by Tassiulas and Ephremides (1992) in the context of communication networks. After that, the max-weight algorithm and the back-pressure algorithm, which is a generalization of the max-weight algorithm, were studied intensively in the literature. The book by Srikant and Ying (2014) provides an excellent summary. The performance of the max-weight algorithm in the context of a switch operating in heavy traffic has been studied by Maguluri and Srikant (2016). The back-pressure algorithm was used in the context of online ad matching in Tan and Srikant (2012) and in the context of ride hailing in Kanoria and Qian (2019).

Heavy traffic analysis of the max-weight algorithm in the context of a single-sided queue has a long line of literature. One analysis approach is based on fluid limits, diffusion limits, and reflected Brownian motion (RBM) (Harrison 2013). In this approach, the queueing process is studied under an appropriate scaling, and the corresponding limiting fluid or diffusion process is shown to converge to a lower-dimensional RBM. This phenomenon is called state space collapse (SSC). If the RBM is single dimensional, then it is called complete resource pooling (CRP). Examples on this line of work to study SSC under the max-weight algorithm in the context of single-sided queues are Williams (1998), Stolyar (2004), and Gamarnik and Zeevi (2006). In this paper, we employ another approach based on the Lyapunov drift method developed by Eryilmaz and Srikant (2012) and later used by Maguluri and Srikant (2015) for switch systems. We generalize the Lyapunov function for two-sided queues and analyze the maxweight algorithm under the CRP condition similar to that in Gurvich and Whitt (2009) and Shi et al. (2019).

1.3. Notation

Throughout the paper, vectors are denoted by boldface letters. Functions applied on vectors are defined entrywise; for example, $F(\lambda)$ is defined to be $(F(\lambda_1), \ldots, \beta_n)$ $F(\lambda_m)$). For any two vectors $\mathbf{a} \in \mathbb{R}^n$ and $\mathbf{b} \in \mathbb{R}^m$, we denote the concatenated vector of dimension n + m by (a,b). We denote the *n*-dimensional vector with all ones by $\mathbf{1}_n$ and the *n*-dimensional vector with all zeroes by $\mathbf{0}_n$; we omit the subscript n if the sizes of these vectors are clear from the context. If x and y are of the same dimension, we use $\langle x,y \rangle$ to denote the inner product and $\mathbf{x} \circ \mathbf{y}$ to denote the Hadamard product (i.e., entrywise product). Any inequality $x \le y$ is also defined entrywise. We use the superscript "s" to denote variables related to servers and the superscript "c" for variables related to customers. We use $\mathbf{e}_{i}^{(c)}$ and $\mathbf{e}_{i}^{(s)}$ to represent unit vectors with a one for type *j* customer and type *i* server, respectively, and all zeroes otherwise.

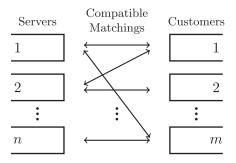
2. Model

We represent the types of customers and servers by a bipartite graph $G(N \cup M, E)$, where N is the set of server types with |N| = n, M is the set of customer types with |M| = m, and E is the set of edges representing customer and server types that are compatible with each other (see Figure 1). A pair $(i,j) \in E$ if and only if a type j customer can be served by a type i server. Each node in the bipartite graph is a queue of customers or servers waiting to be matched with any one of the compatible counterparts. Our convention is to refer to incoming customers as demand and incoming servers as supply.

At each point in time, the system operator posts a price for each customer and server type. Customers willing to pay the quoted prices, as well as servers who are willing to provide their service at the posted prices (i.e., wages), are admitted to the system. Thus, the system operator can vary the prices to control the arrival rates of customers and servers. Customers and servers then wait in queues until they are matched. The FCFS discipline is employed for each queue separately, but it may not hold among different types of customers and servers. After a customer is matched with a compatible server, we assume that they depart from the system instantaneously to complete the service process. The system operator's objective is to find a joint pricing and matching policy under which the system is stable (positive recurrent) and the long-run average profit is maximized.

We assume that customers and servers arrive according to nonhomogeneous Poisson processes. For each server type $i \in N$, there exists a supply curve $\mu_i : \mathbb{R}_+ \to \mathbb{R}_+$, such that if the system operator sets a price $p_i^{(s)}$ and the expected waiting time is $w_i^{(s)}$, the resulting arrival rate is $\mu_i \Big(p_i^{(s)} - s_i^{(s)} w_i^{(s)} \Big)$, where the constant $s_i^{(s)}$ is the unit waiting cost of server type i. Similarly, for each customer type $j \in M$, there exists a demand curve $\lambda_j : \mathbb{R}_+ \to \mathbb{R}_+$, such that if the system operator sets a price $p_i^{(c)}$ and the expected waiting time

Figure 1. Bipartite Graph Representation for Two-Sided Queues



is $w_j^{(c)}$, the resulting arrival rate is $\lambda_j \Big(p_j^{(c)} + s_j^{(c)} w_j^{(c)} \Big)$, where $s_j^{(c)}$ is the unit waiting cost of customer type j. We make the following assumption on the supply and demand curves.

Assumption 1. The supply curves $\mu_i : \mathbb{R}_+ \to \mathbb{R}_+ \ (\forall i \in N)$ are strictly increasing and twice continuously differentiable. The demand curves $\lambda_j : \mathbb{R}_+ \to \mathbb{R}_+ \ (\forall j \in M)$ are strictly decreasing and twice continuously differentiable.

Because λ_j and μ_i are strictly monotone, their inverse functions exist, and we denote them by F_j : $\mathbb{R}_+ \to \mathbb{R}_+$ ($\forall j \in M$) and $G_i : \mathbb{R}_+ \to \mathbb{R}_+$ ($\forall i \in N$), respectively. In addition, we define the revenue and cost functions as $r_j^{(c)}(\lambda_j) \triangleq \lambda_j F_j(\lambda_j)$ for all $j \in M$ and $r_i^{(s)}(\mu_i) \triangleq \mu_i G_i(\mu_i)$ for all $i \in N$. We make the following assumption on the revenue and cost functions.

Assumption 2. The revenue function $r_j^{(c)}(\lambda_j)$ is concave $(\forall j \in M)$. The cost function $r_i^{(s)}(\mu_i)$ is convex $(\forall i \in N)$.

The concavity assumption on revenue function follows from the economic law of diminishing marginal return; as the system operator increases the customer arrival rate λ_j , the marginal revenue $dr_j^{(c)}(\lambda_j)/d\lambda_j$ decreases, which implies that the revenue function $r_j^{(c)}(\lambda_j)$ is concave. This assumption is often assumed in the revenue management literature (Gallego and Van Ryzin 1994, Kim and Randhawa 2017). We assume that the marginal cost $dr_i^{(s)}(\mu_i)/d\mu_i$ increases with μ_i because it becomes harder to recruit servers when we try to increase server arrival rate. This implies that the cost function $r_i^{(s)}$ is convex.

For those customers and servers waiting in queues, the system operator uses matching controls to govern the queueing process. At any given time, suppose $q_i^{(s)}$ is the number of type i servers waiting in queue and

 $q_j^{(c)}$ is the number of type j customers waiting in queue. We denote the vector of all queue lengths by $\mathbf{q} = (q_j^{(c)}, \ \forall j \in M, \ q_i^{(s)}, \ \forall i \in N)$. We denote the number of type i servers to be matched to type j customers by y_{ij} . The set of feasible matching decisions is

$$Y(\mathbf{q}) \triangleq \left\{ \mathbf{y} \in \mathbb{Z}_{+}^{nm} \middle| \sum_{i=1}^{n} y_{ij} \leq q_{j}^{(c)} (\forall j \in M), \right.$$
$$\left. \sum_{j=1}^{m} y_{ij} \leq q_{i}^{(s)} (\forall i \in N), y_{ij} = 0 (\forall (i,j) \notin E) \right\}.$$

We also define the projection of $Y(\mathbf{q})$ to the queuelength space as

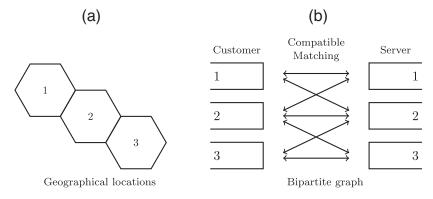
$$X(\mathbf{q}) \triangleq \left\{ \mathbf{x} \in \mathbb{Z}_{+}^{n+m} \middle| \exists \mathbf{y} \in Y(\mathbf{q}) : x_{i}^{(s)} = \sum_{j=1}^{m} y_{ij} (\forall j \in M), \right.$$
$$\left. x_{j}^{(c)} = \sum_{i=1}^{n} y_{ij} (\forall i \in N) \right\}.$$
(1)

When a pair of customer and server is matched by the system, they both depart from the system. Because a customer is only compatible to a subset of server types, the system operator may have an incentive to hold some customers or servers in queue in order to achieve better matches in future.

2.1. Example: Ride Hailing

An application of the two-sided queueing model is in ride-hailing systems. In such a system, the customer and server (drivers) types, as well as the matching compatibility graph, are determined by their geographical locations. A simple example with three regions is shown in Figure 2. (Here, we ignore issues such as vehicle capacity and number of seats requested by customers, which can be accounted for by creating additional customer and server types.) Based on the price and the waiting time quoted to customers, only a

Figure 2. A Ride-Hailing System with Three Regions



Notes. (a) Geographical locations. (b) Bipartite graph. We assume that riders can only be matched to cars in their own region or any neighboring regions. The two-sided system generated from the map is shown in panel (b).

fraction of them who open the app will book a ride, which determines the customer arrival rate. Similarly, based on the price quoted to the drivers, they will choose whether to provide service. Thus, the arrival rates of customer and drivers depend on price and wait time and are governed by the demand and supply curve of each region. After a customer confirms the price and books a ride, the system operator can determine which driver (from what region) should be matched to the customer. If a driver accepts the ride request, then they immediately become unavailable for any other ride requests (departing from the system). After the ride is complete, the car becomes available again, possibly in a different region. A simplifying assumption in our model is that we treat a driver who completes the service and reenters the system the same as a new arrival.

2.2. Continuous Time MDP Formulation

We now formulate the system operator's decision problem as a continuous time Markov decision process (CTMDP) and specify its states, actions, transition rates, and rewards. The system state is represented by the queue lengths of all customer and server types $\mathbf{q} \in \mathbb{Z}_+^{n+m}$. The actions of the CTMDP include both pricing and matching decisions. The matching decision must satisfy $\mathbf{x} \in X(\mathbf{q})$ defined by Equation (1). For the pricing decision, in order to leverage Assumption 2, it is more convenient to use arrival rates (λ, μ) rather than prices as the control variables. In particular, for customer type $j \in N$, setting the arrival rate to λ_j is equivalent to setting the price to $p_i^{(c)} = F_i(\lambda_i) - s_i^{(c)} w_i^{(c)}(\mathbf{q})$. Similarly, for server type $i \in M$, setting the arrival rate to μ_i is equivalent to setting the price to $p_i^{(s)} = G_i^{(s)}(\mu_i) + S_i^{(s)}w_i^{(s)}(\mathbf{q})$. Thus, the action is a tuple $\mathbf{z} \triangleq (\lambda, \mu, x) \in \mathbb{R}^{2(m+n)}$. Given this action, the transition rate from state q to state q + $\mathbf{e}_{i}^{(c)} - \mathbf{x}$ (i.e., having a new arrival of type j customer) is λ_j ($\forall j \in M$), and a reward of $p_j^{(c)}$ is received upon the new arrival. The transition rate from state \mathbf{q} to state $\mathbf{q} + \mathbf{e}_i^{(s)} - \mathbf{x}$ (i.e., having a new arrival of type *i* server) is μ_i ($\forall i \in N$), and a cost of $p_i^{(s)}$ is paid upon the new arrival. The system operator's objective is to find a pricing and matching policy such that the long-run average profit earned by the system operator is maximized. We restrict our attention to policies that make the system stable in the long run, which is defined as follows.

Definition 1. A joint pricing and matching policy is said to be stable if the continuous time Markov chain induced by this policy has a positive recurrent communicating class that contains the state $\mathbf{q} = \mathbf{0}$.

Remark 1 (Average Waiting Time). It is technically challenging to analyze the exact waiting time $w_i^{(s)}(\mathbf{q})$ and $w_j^{(c)}(\mathbf{q})$ because the waiting time of one type may depend on the queue lengths of all the types as well as the policy and matching policy used by the system operator. Additionally, in some applications, real-time queue-length information may not be visible to all market participants (Zohar et al. 2002). Therefore, we make a simplifying assumption that the waiting time *perceived* by the customers and servers is the long-run average waiting time. That is, we assume

$$p_{j}^{(c)} = F_{j}(\lambda_{j}) - s_{j}^{(c)} \mathbb{E}[w_{j}^{(c)}(\mathbf{q})] \quad \forall j \in M,$$

$$p_{i}^{(s)} = G_{i}^{(s)}(\mu_{i}) + s_{i}^{(s)} \mathbb{E}[w_{i}^{(s)}(\mathbf{q})] \quad \forall i \in N.$$

The scheme of announcing the long-run average waiting time to (impatient) customers is commonly assumed in the literature (Zohar et al. 2002, Armony et al. 2009). Additionally, in the large-scale setting that will be considered in the following sections, approximating real-time estimated waiting time with the long-run average waiting time will only result in a negligible error term of a higher order (see Kim and Randhawa 2017, section 6.1 for a similar argument).

2.2.1. Equivalence to Holding Cost Models. The model assumes that customers and servers are sensitive to both prices and waiting costs when they decide to enter the queueing system. We now consider an alternative model, where customers and servers only react to prices, whereas the system operator pays *additional* holding costs for market participants waiting in queues. In particular, in this alternative model, the states, actions, and transition rates remain the same. Given a state \mathbf{q} and an action $\mathbf{z} = (\lambda, \mu, \mathbf{x})$, the reward function is defined as

$$\mathcal{R}(\mathbf{q}, \mathbf{z}) \triangleq \sum_{j=1}^{m} \lambda_{j} F_{j}(\lambda_{j}) - \sum_{i=1}^{n} \mu_{i} G_{i}(\mu_{i}) - \sum_{j=1}^{m} s_{j}^{(c)} q_{j}^{(c)} - \sum_{i=1}^{n} s_{i}^{(s)} q_{i}^{(s)},$$
(2)

where $s_j^{(c)}$ and $s_i^{(s)}$ are the customers' and servers' impatience parameters, respectively, introduced in the original model. The following result shows that the two modeling approaches are indeed equivalent.

Proposition 1. For any given control policy, the delaysensitive model and the holding cost model have the same long-run average profit.

The proof of Proposition 1 follows an application of Little's Law and can be found in Online Appendix A.1. The advantage of considering the holding cost model is that the reward function $\mathcal{R}(\mathbf{q},\mathbf{z})$ does not explicitly depend on the waiting time. Hence, we use the holding cost model in the rest of the paper.

2.3. Discrete Time MDP Formulation by Uniformization

Instead of analyzing the CTMDP directly, we use the well-known uniformization technique (e.g., Puterman 1994, chapter 11) to obtain an equivalent discrete time Markov decision process (DTMDP), which will simplify our analysis. The uniformized process works as follows. We first choose a uniformization parameter \boldsymbol{c} defined here.

Definition 2. Suppose there exist constants $\lambda_{\max} \in \mathbb{R}_+^m$ and $\mu_{\max} \in \mathbb{R}_+^n$ such that for any price vector $\mathbf{p} \ge 0$ we have, $\lambda(\mathbf{p}) \le \lambda_{\max}$ and $\mu(\mathbf{p}) \le \mu_{\max}$. Let c be any constant such that $c \ge \langle \mathbf{1}_m, \lambda_{\max} \rangle + \langle \mathbf{1}_n, \mu_{\max} \rangle$.

The uniformized DTMDP is endowed with the same state \mathbf{q} and action $\mathbf{z} = (\lambda, \mu, \mathbf{x})$ as the CTMDP. Let $Z(\mathbf{q}) = [0, \lambda_{\max}] \cup [0, \mu_{\max}] \cup X(\mathbf{q})$ be the set of feasible actions for queue length $\mathbf{q} \in \mathbb{R}^{m+n}_+$. In the uniformized DTMDP, there is at most one customer arrival or one server arrival in each period. The state transitions from \mathbf{q} to $\mathbf{q} + \mathbf{e}_j^{(c)} - \mathbf{x}$ with probability λ_j/c $(\forall j \in M)$; it transitions from \mathbf{q} to $\mathbf{q} + \mathbf{e}_i^{(s)} - \mathbf{x}$ with probability μ_j/c $(\forall i \in N)$. Otherwise, no arrival happens in this period, and the state remains at \mathbf{q} with probability $1 - ((\mathbf{1}_m, \lambda) + \langle \mathbf{1}_n, \mu \rangle)/c$. The expected reward in one period is given by $\mathcal{R}(\mathbf{q}, \mathbf{z})/c$. Let \mathbf{q}' be the state in the next period. The Bellman equation of the DTMDP is

$$h(\mathbf{q}) + \frac{\gamma}{c} = \max_{\mathbf{z} \in Z(\mathbf{q})} \left\{ \frac{\mathcal{R}(\mathbf{q}, \mathbf{z})}{c} + \mathbb{E}[h(\mathbf{q}') \mid \mathbf{q}, \mathbf{z}] \right\}, \quad \forall \mathbf{q} \in \mathbb{Z}_{+}^{n+m},$$
(3)

where

$$\mathbb{E}[h(\mathbf{q}') \mid \mathbf{q}, \mathbf{z}] = \sum_{j=1}^{m} \frac{\lambda_j}{c} h(\mathbf{q} + \mathbf{e}_j^{(c)} - \mathbf{x}) + \sum_{i=1}^{n} \frac{\mu_i}{c} h(\mathbf{q} + \mathbf{e}_i^{(s)} - \mathbf{x}) + \left(1 - \sum_{j=1}^{m} \frac{\lambda_j}{c} - \sum_{i=1}^{n} \frac{\mu_i}{c} h(\mathbf{q})\right).$$
(4)

In the equation, the solution γ is the optimal long-run average profit, and $h(\mathbf{q})$ is the bias function associated with state \mathbf{q} ($\forall \mathbf{q} \geq 0$). (Note that the optimal solution of the uniformized DTMDP satisfies the Bellman equation because we require the optimal policy to be stable (see Definition 1).) The term $\mathbb{E}[h(\mathbf{q}') \mid \mathbf{q}, \mathbf{z}]$ is the expectation of the bias function h after one transition in the uniformized process. The expectation is taken with respect to the one-period transition probabilities conditional on the state \mathbf{q} and the action \mathbf{z} .

In Online Appendix A, we present additional analysis of the uniformized DTMDP. We show the monotonicity structure of the optimal pricing policy in the single-link queueing system (i.e., m = n = 1). Unfortunately, as the number of customer and server types becomes large, solving the DTMDP becomes intractable because of the curse of dimensionality. We propose two approximation methods to obtain near-optimal solutions to the DTMDP.

The first method is based on fluid approximation. The remainder of the paper primarily focuses on this approach. The second method uses value function approximation. We defer details of the second method to Online Appendix A, as the remaining parts of the paper do not rely on it.

2.4. Max-Weight Matching Policy

In the following sections, we will extensively use the *max-weight* matching policy, so we provide its definition here. Suppose the system has state \mathbf{q} , and the set of feasible matches is $X(\mathbf{q})$ (see Equation (1)). The policy chooses the matching decision \mathbf{x} to be the solution of

$$\underset{\mathbf{x} \in X(\mathbf{q})}{\text{arg max}} \left\{ \langle \mathbf{q}, \mathbf{x} \rangle \right\}. \tag{5}$$

In other words, under the max-weight policy, when there is either a customer or a server arrival, a match will be made if any of the compatible types has a non-empty queue, and we will always match the arriving customer/server to the compatible type with the most customers/servers waiting in queue. Otherwise, if all the compatible counterparts' queues are empty, then the arrival is inserted into the queue of its own type.

The max-weight matching policy, originally proposed by Tassiulas and Ephremides (1992), is extensively studied in the queueing literature. This literature is reviewed in Section 1.2. Apart from the queueing literature, in our model specifically, there is also an alternative way to motivate the max-weight matching policy through the quadratic value function approximation of the MDP. Suppose the bias function in Equation (4) is approximated by $h(\mathbf{q}) \approx \langle 1, \mathbf{q}^2 \rangle$; then, the *optimal* matching policy of the DTMDP will be very close to the max-weight policy defined in Equation (5). Online Appendix A.3 contains a detailed discussion of the value function approximation method.

```
Algorithm 1 (Max-Weight Matching Policy)
```

output: matching decision $\mathbf{y}(k)$

```
input: current queue length \mathbf{q}(k), new arrival \mathbf{a}(k) {k
is a decision epoch
initialization: y(k) = 0
for i \in N do
   if a_i^{(s)}(k) = 1 and \max_{j:(i,j)\in E} q_i^{(c)} > 0 then
     choose j^* \in \arg\max_{j:(i,j)\in E} q_j^{(c)}
                                                (breaking
      arbitrarily)
      set y_{ij^*}(k) = 1
   end if
end for
for j \in M do
   if a_i^{(c)}(k) = 1 and \max_{i:(i,j)\in E} q_i^{(s)} > 0 then
      let i^* \in \arg\max_{i:(i,j)\in E} q_i^{(s)} (breaking ties arbitrarily)
      set y_{i*i}(k) = 1
   end if
end for
```

3. Asymptotically Optimal Policies in the Large-Scale Regime

3.1. Fluid Model and Large-Scale Regime

We consider a fluid approximation of the queueing system where random arrivals are replaced by deterministic arrival processes. The fluid model is a deterministic optimization problem maximizing the long-run average profit. Suppose customers arrive with constant rates λ and servers arrive with constant rates μ . Let χ_{ij} be the average rate of type i server matched to the type j customer for all $(i,j) \in E$. The fluid model is defined as

$$\gamma^* = \max_{\lambda, \mu, \chi} \langle F(\lambda), \lambda \rangle - \langle G(\mu), \mu \rangle$$
 (6a)

subject to
$$\lambda_j = \sum_{i=1}^n \chi_{ij}, \quad \forall j \in M,$$
 (6b)

$$\mu_i = \sum_{j=1}^m \chi_{ij}, \quad \forall i \in N, \tag{6c}$$

$$\chi_{ij} = 0, \ \forall (i,j) \notin E, \ \chi_{ij} \ge 0, \ \forall (i,j) \in E.$$
 (6d)

We denote an optimal solution to the fluid problem by $(\lambda^*, \mu^*, \chi^*)$.

To interpret the fluid model, note that Equations (6b) and (6c) are the balance equations for the number of customers and servers matched. Equation (6d) specifies that matching is only allowed among compatible customer-server pairs. Intuitively, it is easy to see that these constraints are necessary because if the balance equations do not hold, then some customer or server types will keep accumulating over time. Thus, the optimization program Equation (6) serves as an *upper bound* on the achievable profit under any pricing and matching policy that makes the system stable. This is formally shown in the following proposition. The proof can be found in Online Appendix B.1.

Proposition 2. The optimal value of the fluid problem Equation (6) is an upper bound on the long-run expected profit under any policy that makes the system stable.

In the remainder of this section, we analyze the two-sided queueing system in a large-scale regime where the arrival rates of all customer and server types are simultaneously scaled by a factor of $\eta \in \mathbb{N}$.

Definition 3 (Large-Scale Regime). Consider a family of two-sided queueing systems associated with the same bipartite graph $G(N \cup M, E)$ parametrized by $\eta \in \mathbb{N}$. For the η th system, the demand and supply curves satisfy $F^{\eta}(\eta \lambda) = F(\lambda)$ for all $\mathbf{0}_m \leq \lambda \leq \lambda_{\max}$ and $G^{\eta}(\eta \mu) = G(\mu)$ for all $\mathbf{0}_n \leq \mu \leq \mu_{\max}$.

The large-scale regime is commonly assumed in the dynamic pricing and matching literature (Gurvich and

Ward 2014, Özkan and Ward 2020), which ensures that supply and demand are balanced as the system scales up. According to Definition 3, it is easily verified that the fluid solution to the η th scaled system is given by $\eta \lambda^*$ and $\eta \mu^*$, where λ^* and μ^* are the optimal solution of the unscaled fluid model Equation (6).

Definition 4 (Profit Loss). The profit loss (denoted by L^{η}) of a policy is the difference between the optimal value of the (scaled) fluid model, denoted by γ_*^{η} , and the long-run average profit (including the penalty incurred because of waiting) under that policy.

The optimal value of the η th fluid model is $\gamma^{\eta}_{*} = \eta \gamma^{*}$. Therefore, if the profit loss of a policy is sublinear in η , namely $L^{\eta} = o(\eta)$, we say the policy is asymptotically optimal in the large-scale regime.

3.2. Fluid Pricing Policy

Based on the fluid model, we propose a static pricing policy defined as follows:

$$\lambda_{j}(\mathbf{q}) = \begin{cases} \lambda_{j}^{*} & \text{if } q_{j}^{(c)} < q_{\max}^{\eta} \\ 0 & \text{otherwise} \end{cases} \quad \forall j \in M,$$

$$\mu_{i}(\mathbf{q}) = \begin{cases} \mu_{i}^{*} & \text{if } q_{i}^{(s)} < q_{\max}^{\eta} \\ 0 & \text{otherwise} \end{cases} \quad \forall i \in N.$$

$$(7)$$

Here, q_{max}^{η} denotes the maximum queue buffer size; it is a parameter that depends on η , which will be specified later.

The main intuition of the fluid pricing policy is the following. When all queues are below their maximum buffer capacity \mathbf{q}^η , the profit rate of the fluid pricing policy is exactly equal to $\eta \gamma^*$. If any customer queue is full, say $q_j^{(c)} = q_{\max}^\eta$, then all future arrivals to queue j will be rejected until at least one customer waiting in queue j is matched. Thus, a fraction of revenue is lost because of customer rejections. More specially, let γ^η be the long-run average profit of the fluid pricing policy (excluding waiting costs). Let $\mathbf{I}^{(s)}(q_{\max}^\eta)$ be a (vector) indicator function representing whether server queues are at the maximum capacity, and let $\mathbf{I}^{(c)}(q_{\max}^\eta)$ be a (vector) indicator function representing whether customer queues are at the maximum capacity. Then, we have

$$L^{\eta} = \gamma_{*}^{\eta} - (\gamma^{\eta} - \langle \mathbf{s}, \mathbb{E}[\mathbf{q}] \rangle)$$

$$= \eta(\langle F(\boldsymbol{\lambda}^{*}), \boldsymbol{\lambda}^{*} \rangle - \langle G(\boldsymbol{\mu}^{*}), \boldsymbol{\mu}^{*} \rangle)$$

$$- \langle F(\boldsymbol{\lambda}^{*}), \eta \boldsymbol{\lambda}^{*} \circ (\mathbf{1} - \mathbb{E}[\mathbf{I}^{(c)}(q_{\max}^{\eta})]) \rangle$$

$$- \eta \langle G(\boldsymbol{\mu}^{*}), \boldsymbol{\mu}^{*} \circ (\mathbf{1} - \mathbb{E}[\mathbf{I}^{(s)}(q_{\max}^{\eta})]) \rangle + \langle \mathbf{s}, \mathbb{E}[\mathbf{q}] \rangle$$

$$= \eta \Big(\langle F(\boldsymbol{\lambda}^{*}), (\boldsymbol{\lambda}^{*} \circ \mathbb{E}[\mathbf{I}^{(c)}(q_{\max}^{\eta})]) \rangle$$

$$- \langle G(\boldsymbol{\mu}^{*}), (\boldsymbol{\mu}^{*} \circ \mathbb{E}[\mathbf{I}^{(s)}(q_{\max}^{\eta})] \rangle \Big) + \langle \mathbf{s}, \mathbb{E}[\mathbf{q}] \rangle, \tag{8}$$

where the first equality follows from Definition 4 and the second equality uses the definition of the fluid pricing policy. As a result, Equation (8) shows that the profit loss of the fluid pricing policy depends on the parameter q_{\max}^{η} . If we increase the buffer capacity q_{\max}^{η} , then the probability of dropping customers/servers will reduce (i.e., $\mathbb{E}[\mathbf{I}(\mathbf{q}_{\max}^{\eta})]$ will decrease). However, increasing the buffer capacity will lead to increasing in the expected queue lengths, which will increase the penalty incurred because of waiting. Thus, we choose buffer capacity to balance the tradeoff in order to minimize the overall profit loss. Precisely, we will see that choosing $q_{\max}^{\eta} \sim \sqrt{\eta}$ will result in $\mathbb{E}[\mathbf{I}(\mathbf{q}_{\max}^{\eta})] \sim \eta^{-1/2}$ and $\mathbb{E}[\langle \mathbf{1}_{m+n}, \mathbf{q} \rangle] \sim \sqrt{\eta}$, which attains the optimal profit loss.

Theorem 1. Suppose a family of two-sided queues is given by the bipartite graph $G(N \cup M, E)$ parameterized by η . The profit loss L^{η} under the fluid pricing policy Equation (7) and max-weight matching (Algorithm 1) is $O(\sqrt{\eta})$, where $q_{\max}^{\eta} = \gamma \sqrt{\eta}$ for any positive constant γ .

The proof of Theorem 1 can be found in Online Appendix B.2. In addition, it can be shown that the $O(\sqrt{\eta})$ profit loss rate cannot be improved using any fluid pricing policy. The proof of the proposition is presented in Online Appendix B.3.

Proposition 3. For a family of single link two-sided queues parametrized by η , any fluid pricing policy will have a profit loss L^{η} that is at least $\Omega(\sqrt{\eta})$. The choice of $q_{\max}^{\eta} = \gamma \sqrt{\eta}$ for any positive constant γ provides the optimal profit loss rate $\Theta(\sqrt{\eta})$.

3.3. Two-Price Policy

A main drawback of the fluid pricing policy is that the prices are not adaptive to changes in the system state. In this section, we consider another policy that uses *two* different prices for each customer/server type. The proposed two-price policy is built on the two-price policy in Kim and Randhawa (2017) for single-server queues. Our contribution lies in a joint analysis of two-price and dynamic matching policies in a multitype queueing network.

The two-price policy can be viewed as a generalization of the fluid pricing policy. We introduce additional parameters $\boldsymbol{\theta} \in \mathbb{R}^m_+$, $\boldsymbol{\phi} \in \mathbb{R}^n_+$, and $\sigma^\eta > 0$, which govern the arrival rates of the customers and servers, respectively, when the queue length is greater than a certain threshold τ^η_{max} . The two-price policy is defined as

$$\lambda_{j}(\mathbf{q}) = \begin{cases} \eta \lambda_{j}^{*} & \text{if } q_{j}^{(c)} \leq \tau_{\max}^{\eta} \\ \eta \lambda_{j}^{*} - \theta_{j} \sigma^{\eta} & \text{otherwise} \end{cases} \quad \forall j \in M,$$

$$\mu_{i}(\mathbf{q}) = \begin{cases} \eta \mu_{i}^{*} & \text{if } q_{i}^{(s)} \leq \tau_{\max}^{\eta} \\ \eta \mu_{i}^{*} - \phi_{i} \sigma^{\eta} & \text{otherwise} \end{cases} \quad \forall i \in N.$$

$$(9)$$

The policy sets a threshold τ_{\max}^{η} for all customer and server types. It uses the fluid arrival rates when queue lengths are below this threshold and then reduces the arrival rates by $\theta_i \sigma^{\eta}$ outside this threshold for type *j* customer. Similarly, the policy reduces the server arrival rates outside the threshold by $\phi_i \sigma^{\eta}$ for type *i* server. Here, τ_{\max}^{η} , σ^{η} , $\boldsymbol{\theta}$, and $\boldsymbol{\phi}$ are parameters that will be specified later. (Our convention is to use superscript η to denote any parameter or quantity that is associated with the η th scaled system.) Intuitively, for any type of customer/server, if we increase σ^{η} , the queue length will have a larger negative drift when it exceeds the threshold τ_{\max}^{η} , so the expected queue length $\mathbb{E}[\langle \mathbf{1}_{m+n}, \mathbf{q} \rangle]$ will be smaller. However, if σ^{η} are too large, the arrival rates outside the threshold au_{max}^η will be far from the optimal fluid arrival rates, which will result in a larger profit loss. Thus, there is a tradeoff between the expected queue length and profit loss. For the matching algorithm associated with the twoprice policy, here we use the max-weight matching algorithm as defined in Equation (5). (Other matching algorithms will be considered in Section 4.2.) The following theorem provides a bound on the asymptotic performance of the two-price policy as η tends to

Theorem 2. Consider a family of two-sided queues parametrized by η represented by the bipartite graph $G(N \cup M, E)$. The profit loss L^{η} under the two-price policy Equation (9) and the max-weight matching (Algorithm 1) is $O(\eta^{1/3})$ for any $\tau_{\max}^{\eta} \leq \eta^{1/3}$, $\sigma^{\eta} = \eta^{2/3}$ and constants $\theta > \mathbf{0}_m$, $\phi > \mathbf{0}_n$.

The theorem shows that the profit loss of the two-price policy is $O(\eta^{1/3})$, which is better than the $O(\sqrt{\eta})$ loss in the fluid pricing policy. The proof of the theorem contains two main steps. The first step is to show that the system is stable under the two-price policy and that the expected queue lengths are bounded. We also give an upper bound of the expected queue lengths (Lemma 1). The second step in the proof is to estimate the profit loss L^{η} (Lemma 2) by applying the Karush-Kuhn-Tucker (KKT) conditions of the fluid problem.

Lemma 1. For a system of two-sided queues operating under the two-price policy and the max-weight matching algorithm parameterized by η , the system is positive recurrent for any $\theta > 0_m$, $\phi > 0_n$, $\sigma^{\eta} > 0$, and $\tau_{\max}^{\eta} > 0$. The expected queue lengths are bounded by

$$\begin{split} & \mathbb{E}\left[\left\langle \boldsymbol{\theta}, \mathbf{q}^{(c)} \right\rangle\right] + \mathbb{E}\left[\left\langle \boldsymbol{\phi}, \mathbf{q}^{(s)} \right\rangle\right] \\ & \leq \tau_{\max}^{\eta} \left(\sum_{j=1}^{m} \theta_{j} \mathbb{P}\left[q_{j}^{(c)} > \tau_{\max}^{\eta}\right] + \sum_{i=1}^{n} \phi_{i} \mathbb{P}\left[q_{i}^{(s)} > \tau_{\max}^{\eta}\right]\right) \\ & + \frac{\eta}{\sigma^{\eta}} \left(\left\langle \mathbf{1}_{n}, \boldsymbol{\mu}^{*} \right\rangle + \left\langle \mathbf{1}_{m}, \boldsymbol{\lambda}^{*} \right\rangle\right). \end{split}$$

Lemma 2. For a system of two-sided queues operating under the two-price policy and the max-weight matching policy, for any $\theta > 0_m$, $\phi > 0_n$, and $\tau_{\max^{\eta}} > 0$, we have

$$\begin{split} &\sum_{j \in M} \left(F_j'(\lambda_j^*) \lambda_j^* + F_j(\lambda_j^*) \right) \theta_j \mathbb{P}[q_j^{(c)} > \tau_{\max}^{\eta}] \\ &- \sum_{i \in N} \left(G_i'(\mu_i^*) \mu_i^* + G_i(\mu_i^*) \right) \phi_i \mathbb{P}[q_i^{(s)} > \tau_{\max}^{\eta}] \\ &\leq & |E|^2 \max_{i \in N, i \in M} \left\{ \phi_i, \theta_j \right\} \frac{\sigma^{\eta}}{n}. \end{split}$$

3.4. Lower Bound

In this section, we will obtain a lower bound on the profit loss under a broad family of policies and thus, establish that the $O(\eta^{1/3})$ rate obtained by the two-price policy in Theorem 2 is optimal. In particular, we consider a family of pricing policies that have the following form:

$$\lambda_{j} = \eta \lambda_{j}^{*} + f_{j} \left(\frac{\mathbf{q}}{n^{\alpha}} \right) \eta^{\beta} \quad \forall j \in M, \tag{10}$$

$$\mu_i = \eta \mu_i^* + g_i \left(\frac{\mathbf{q}}{\eta^{\alpha}} \right) \eta^{\beta} \quad \forall i \in \mathbb{N}.$$
 (11)

The motivation for this policy is as follows. The first terms in Equations (10) and (11) (i.e., $\eta \lambda_j^*$ and $\eta \mu_i^*$, respectively) are static and result from the solution of the fluid model; the second terms account for dynamic adjustments as the queue length changes. We assume the adjustment terms can be further decomposed into two terms: a function that rescales the queue length, $f_j(\cdot)$ or $g_i(\cdot)$, and a term that determines the scaling of price adjustments, η^β , for some $1 > \beta > 0$. As the arrival rates are scaled up, the average queue length will also increase. Thus, we rescale the queue length in the functions $f_j(\cdot)$ and $g_i(\cdot)$ for all $i \in N$ and $j \in M$ by η^α for some $1 \ge \alpha \ge 0$.

For our analysis, we assume the pricing policy to satisfy the following conditions.

Condition 1.

a. There exist constants $\Gamma \in \mathbb{R}_+^m$ and $\Psi \in \mathbb{R}_+^n$ such that $|f_j(\mathbf{q}/\eta^\alpha)| \le \Gamma_j$ for all $j \in M$ and $|g_i(\mathbf{q}/\eta^\alpha)| \le \Psi_i$ for all $i \in N$ for all $\mathbf{q} \in S$ and for all $\eta \ge 1$.

b. $0 < \alpha + \beta \le 1$.

c. There exist constants $\kappa > 0$ and $\delta > 0$ such that for all $j \in M$, if $q_j^{(c)}/\eta^{\alpha} > \kappa$, then either $f_j(\mathbf{q}/\eta^{\alpha}) < -\delta$ or there exists $i:(i,j) \in E$ such that $g_i(\mathbf{q}/\eta^{\alpha}) > \delta$ for all η . Similarly, for all $i \in N$, if $q_i^{(s)}/\eta^{\alpha} > \kappa$, then either $g_i(\mathbf{q}/\eta^{\alpha}) < -\delta$ or there exists $j:(i,j) \in E$ such that $f_j(\mathbf{q}/\eta^{\alpha}) > \delta$ for all η .

We now interpret the conditions. Condition 1(a) requires the functions f and g to be bounded given appropriately scaling of the queue lengths \mathbf{q} as η increases. Condition 1(b) states that the rate of queue-length

rescaling (α) should not exceed the rate of rescaling pricing adjustment terms $(1-\beta)$. This condition is needed so that the price adjustment terms are sufficiently large to make the system stable. (In the special case of a single-link system, this assumption is not needed; the extension is presented later in Proposition 4.) Condition 1(c) states that if a queue is too long, we should either decrease the arrival rate of this queue or increase the arrival rates of those matched to this queue.

Aside from the conditions, the pricing forms in Equations (10) and (11) are fairly general because the pricing function of any queue can depend on the entire system state vector (\mathbf{q}) , and we do not make any strong assumptions, such as monotonicity, continuity, or differentiability, on functions f and g. Finally, we emphasize that our analysis does not require any assumption on the form of matching policies.

The two-price policy in Section 3.3 satisfies the condition with

$$f_{j}(\mathbf{q}) = -\theta_{j} \mathbb{1}_{q_{j}^{(c)} > \tau_{\max}} (\forall j \in M),$$

$$g_{i}(\mathbf{q}) = -\phi_{i} \mathbb{1}_{q_{i}^{(s)} > \tau_{\max}} (\forall i \in N), \quad \beta = 2/3.$$
(12)

Now, we present the result on the lower bound.

Theorem 3. For a two-sided queue defined by a graph $G(N \cup M, E)$ operating under any pricing policy of the form Equations (10) and (11) that satisfies Condition 1, if the resulting system is stable, there exists a constant K(F, G, f, g) such that

$$L^{\eta} \ge K \eta^{1/3}.$$

The details of the proof are deferred to Online Appendix D.1. We present an intuitive explanation of the rate in the lower bound.

Remark 2 (Intuitive Explanation of $\eta^{1/3}$). The main reason why the profit loss lower bound is of order $O(\eta^{1/3})$ is because of the trade-off between the expected queue length and the loss in revenue. Consider a pricing policy that deviates from the fluid optimal pricing policy by $\epsilon > 0$; that is, for all $\mathbf{q} \in S$, we have $|\lambda_j(\mathbf{q}) - \lambda_i^*| < \epsilon$ for all $j \in M$ and $|\mu_i(\mathbf{q}) - \mu_i^*| < \epsilon$ for all $i \in N$. One can show that under such a policy, the expected queue length is of the order $1/\epsilon$ and revenue loss is of the order $\eta \epsilon^2$. Specifically, the queue length can be coupled to that of an M/M/1 queue in heavy traffic with parameter ϵ , whose mean queue length is known to be of the order $1/\epsilon$ by the Kingman's bound. The loss in revenue can be estimated by the Taylor series expansion of the revenue function. As the arrival rates under the given pricing policy is close to the optimal fluid arrival rates, the first-order term vanished, and the dominant term is of the second order (namely, $\eta \epsilon^2$). The coefficient is this term is shown to be strictly

positive by analyzing the tail probabilities. Therefore, we have

$$\mathbb{E}[\langle \mathbf{1}_{n+m}, \mathbf{q} \rangle] \sim \frac{1}{\epsilon}, \quad \gamma_*^{\eta} - \gamma^{\eta} \sim \eta \epsilon^2.$$

To achieve the optimal trade-off between expected queue length and profit loss, we choose $\epsilon \sim \eta^{-1/3}$, which results in the $\eta^{1/3}$ profit loss in Theorem 3.

We can further relax Condition 1(b) in the special case of a single-link system (m = n = 1) operating under any two-price policy. The result is stated here, and the proof can be found in Online Appendix D.2.

Proposition 4. For a family of single-link two-sided queues parametrized by η , any two-price policy given by Equation (9) with $\sigma^{\eta} = \eta^{\beta}$ for some $\beta < 1$ and $\tau_{\max}^{\eta} = \eta^{\alpha}$ for some $\alpha \in \mathbb{R}_+$ will have a profit loss L^{η} at least $\Omega(\eta^{1/3})$. The choice of $\tau_{\max}^{\eta} = 1/3$ and $\sigma^{\eta} = \eta^{2/3}$ and any positive constants $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ provides the optimal profit loss $\Theta(\eta^{1/3})$.

Before concluding this section, note that without any further assumptions, the dependence of profit loss L^η on the number of customer types n and the number of server types m can be linear in the worst case. To see this, if the system consisted of n-independent single-link two-sided queues (i.e., the bipartite graph is such that |M| = |N| = n and $E = \{(i, i) : i \in M\}$), then by Theorem 3, the profit loss L^η is trivially lower bounded by $nK\eta^{1/3}$. In the next section, we impose additional conditions on the bipartite graph and show tight lower and upper bounds on the profit loss in terms of the system size.

4. Asymptotically Optimal Policies in the Large-System Regime: The Superiority of Max-Weight Matching

In this section, we present further insights into the max-weight matching algorithm in a large-system regime, in which both the arrival rates and the numbers of customer and server types increase. First, we will show that max weight is delay optimal. Under the fluid pricing policy, max-weight matching minimizes the probability of hitting the queue-length threshold among all matching policies. Under the two-price policy, max-weight matching minimizes the expected sum of queue lengths among all possible matching policies. Second, we compare max-weight matching with a randomized matching policy with probabilities specified by the fluid model and show that max weight has smaller loss in terms of the number of customer/server types. Third, we prove that the profit loss of maxweight matching achieves a tight lower bound in the large-system regime. Together, these results show the superiority of the max-weight policy.

We start by establishing the *state space collapse* under max-weight matching. State space collapse means that

all the customer queues are almost equal in length and that all the server queues are almost equal in length; hence, with high probability, only customers or only servers are waiting in the system. This implies that max weight ends up matching the maximum possible number of customer-server pairs, as only the excess customers/servers are waiting in the system. To achieve the state space collapse, we propose a complete resource pooling condition on the compatibility graph. Similar conditions have been proposed for single-sided queues (Gurvich and Whitt 2009, assumption 2.4; Shi et al. 2019, definition 1).

Condition 2 (CRP). There exists an optimal solution (λ^*, μ^*) to the fluid problem Equation (6) such that for all $J \subseteq M$ and for all $I \subseteq N$, it holds that

$$\sum_{j \in J} \lambda_j^* < \sum_{i: \exists j \in J, (i,j) \in E} \mu_i^*, \quad \sum_{i \in I} \mu_i^* < \sum_{j: \exists i \in I, (i,j) \in E} \lambda_j^*.$$

It is straightforward to verify that the CRP condition implies the connectedness of the graph $G(N \cup M, E)$. The CRP condition also implies that the optimal solution of the fluid problem is in the interior of the feasible region. The following lemma formalizes this observation.

Lemma 3. If Condition 2 is satisfied, there exists $\chi^* \ge 0$ such that $\chi_{ij}^* > 0$ for all $(i,j) \in E$ and $(\lambda^*, \mu^*, \chi^*)$ is an optimal solution to the fluid problem Equation (6).

(All the proofs in this section can be found in Online Appendix E.) The result is not surprising as it is known in the heavy traffic literature (Eryilmaz and Srikant 2012, Lange and Maguluri 2019) that if the arrival rate is approaching a point on the boundary of the capacity region in the interior of a facet, then the system exhibits complete resource pooling.

However, the analysis of state space collapse for two-sided queues does not follow immediately from the literature of single-sided queues and is more involved. We propose a Lyapunov function approach and use the drift method to show state space collapse. To simplify the analysis, in this section we restrict to a setting where m = n, and there exists a perfect matching in the graph $G(N \cup M, E)$.

Condition 3. The graph $G(N \cup M, E)$ has a perfect matching. Without loss of generality, we assume that server type i is connected to customer type i for all $i \in [n]$.

In general, if $m \neq n$ and if the condition is not satisfied, we show in Online Appendix E.6 that the pricing and matching problem under a given general graph can be reformulated as a problem under a new graph where the condition is satisfied. Thus, the results in the following propositions and the theorem can be applied (with minor modifications as shown in Online Appendix E.6) even when the condition does not hold.

4.1. Delay Optimality of Max-Weight Matching

We first show the delay optimality of the max-weight matching algorithm among all possible matching algorithms under the fluid pricing policy.

Proposition 5. *Under the fluid pricing policy with any matching algorithm, we have*

$$q_{\max}^{\eta} \left(\sum_{j=1}^{n} \lambda_{j}^{*} \mathbb{P} \left[q_{j}^{(c)} = q_{\max}^{\eta} \right] + \sum_{i=1}^{n} \mu_{i}^{*} \mathbb{P} \left[q_{i}^{(s)} = q_{\max}^{\eta} \right] \right)$$

$$\geq \frac{\langle \mathbf{1}_{n}, \boldsymbol{\lambda}^{*} \rangle + \langle \mathbf{1}_{n}, \boldsymbol{\mu}^{*} \rangle}{2n + 1/q_{\max}^{\eta}}.$$

Furthermore, under the fluid pricing policy with the maxweight matching algorithm, if $q_{\max}^{\eta} \to \infty$ as $\eta \to \infty$, we have

$$\begin{split} &\lim_{\eta \to \infty} q_{\max}^{\eta} \left(\sum_{j=1}^{n} \lambda_{j}^{*} \mathbb{P} \left[q_{j}^{(c)} = q_{\max}^{\eta} \right] + \sum_{i=1}^{n} \mu_{i}^{*} \mathbb{P} \left[q_{i}^{(s)} = q_{\max}^{\eta} \right] \right) \\ &= \frac{\langle \mathbf{1}_{n}, \boldsymbol{\lambda}^{*} \rangle + \langle \mathbf{1}_{n}, \boldsymbol{\mu}^{*} \rangle}{2n}. \end{split}$$

The proposition states that the max-weight algorithm (asymptotically) minimizes the proportion of time spent in the threshold state among all possible matching algorithms, hence minimizing the revenue loss caused by hitting the queue-length thresholds.

Similarly, the max-weight matching algorithm is delay optimal under the two-price policy. The following proposition states that the max-weight algorithm (asymptotically) minimizes the expected total queue length under the two-price algorithm among all possible matching algorithms.

Proposition 6. Under the two-price policy with $\theta = \phi = \mathbf{1}_n$ and any matching policy, the expected total queue length satisfies

$$\frac{\sigma^{\eta}}{\eta} \mathbb{E}[\langle \mathbf{1}_{2n}, \mathbf{q} \rangle] \geq \frac{\langle \mathbf{1}_{n}, \boldsymbol{\lambda}^{*} \rangle + \langle \mathbf{1}_{n}, \boldsymbol{\mu}^{*} \rangle}{2n}.$$

Furthermore, under the two-price policy with $\theta = \phi = \mathbf{1}_n$ and the max-weight matching policy, if $\lim_{\eta \to \infty} \sigma^{\eta}/\eta = 0$ and $\lim_{\eta \to \infty} \sigma^{\eta} \tau_{\max}^{\eta}/\eta = 0$, we have

$$\lim_{\eta\to\infty}\frac{\sigma^\eta}{\eta}\mathbb{E}[\langle\mathbf{1}_{2n},\mathbf{q}\rangle]=\frac{\langle\mathbf{1}_n,\boldsymbol{\lambda}^*\rangle+\langle\mathbf{1}_n,\boldsymbol{\mu}^*\rangle}{2n}.$$

Notice that the queue-length bound in Proposition 6 is tighter than the bound in Lemma 1 because the former requires the CRP condition (Condition 2), whereas the latter does not require such condition. Together, Propositions 5 and 6 establish the asymptotic delay optimality of the max-weight algorithm.

4.2. Max-Weight Vs. Randomized Matching

In this section, we compare the max-weight policy with a randomized matching policy (defined in Algorithm 2) resulting from the fluid model. The randomized matching algorithm matches an incoming arrival to compatible types at fixed probabilities, which are determined by the fluid solution χ^* (see Equation (6)). If some queues are empty, the probabilities are rescaled proportionally to match only nonempty queues. Unlike the max-weight algorithm, the randomized matching algorithm does not use information about the queue lengths (except for the emptiness of the queues).

Algorithm 2 (Randomized Matching (Nonempty Queues First))

```
input: new arrival \mathbf{a}(k), queue length \mathbf{q}(k), the fluid solution \boldsymbol{\chi}^* {k is a decision epoch} initialization: \mathbf{y}(k) = \mathbf{0} for i \in N do

if a_i^{(s)}(k) = 1 then

set y_{ij}(k) = 1 with probability \frac{\chi_{ij}^*\mathbb{1}_{\{q_j^{(c)}>0\}}}{\sum_{j'=1}^m \chi_{ij'}^*\mathbb{1}_{\{q_j^{(c)}>0\}}} for all j \in M.

end if end for for j \in M do

if a_j^{(c)}(k) = 1 then

set y_{ij}(k) = 1 with probability \frac{\chi_{ij}^*\mathbb{1}_{\{q_i^{(s)}>0\}}}{\sum_{i'=1}^n \chi_{i'j}^*\mathbb{1}_{\{q_i^{(s)}>0\}}} for all i \in N.

end if end for output: matching decision \mathbf{y}(k)
```

We analyze the profit losses of these two matching algorithms and its dependence on the number of customer/server types n when $\eta \to \infty$. First, we consider the fluid pricing policy. The theorem shows that even though both max-weight and randomized matching have $O(\eta^{1/2})$ profit loss, max-weight matching is order $n^{1/2}$ better than randomized matching policy.

Theorem 4. Suppose a family of two-sided queues is given by the bipartite graph $G(N \cup M, E)$ parametrized by η . Under the fluid price policy Equation (7) and randomized matching policy (Algorithm 2), for $q_{\max}^{\eta} = \gamma \eta^{1/2}$, we have $L^{\eta} = O(\eta^{1/2})$. For any $\gamma > 0$, there exists $(\lambda^*, \mu^*, \chi^*)$ satisfying Conditions 2 and 3 such that

$$\liminf_{\eta \to \infty} \frac{L^{\eta}}{\eta^{1/2}} = \Omega(n).$$

In addition, under the fluid price policy (7) and max-weight matching (5) for $q_{\text{max}}^{\eta} = \sqrt{\eta/n}$, we have

$$\limsup_{\eta \to \infty} \frac{L^{\eta}}{\eta^{1/2}} \le n^{1/2} \left(\frac{\langle \mathbf{1}_n, \boldsymbol{\lambda}^* \rangle + \langle \mathbf{1}_n, \boldsymbol{\mu}^* \rangle}{2n} \max_{j \in N} F_j(\lambda_j^*) + 2 \max_{i \in N, j \in M} \left\{ s_i^{(s)}, s_j^{(c)} \right\} \right) = O(n^{1/2}).$$

Next, we compare the max-weight and randomized matching algorithms for the two-price pricing policy.

The theorem shows that both algorithms achieve $O(\eta^{1/3})$ profit loss, whereas max-weight is order $n^{2/3}$ better than randomized matching.

Theorem 5. Suppose a family of two-sided queues is given by the bipartite graph $G(N \cup M, E)$ parametrized by η . Under the two-price policy Equation (9) and randomized matching policy (Algorithm 2), for $\sigma^{\eta} = \eta^{2/3}$ and $\tau^{\eta}_{\max} = \gamma \eta^{1/3}$, we have $L^{\eta} = O(\eta^{1/3})$. For any choice of $\theta > 0$, $\phi > 0$, and $\gamma > 0$, there exists $(\lambda^*, \mu^*, \chi^*)$ satisfying Conditions 2 and 3 such that

$$\liminf_{\eta \to \infty} \frac{L^{\eta}}{\eta^{1/3}} = \Omega(n).$$

In addition, under the two-price policy Equation (9) and max-weight matching Equation (5) with $\theta = \mathbf{1}_n$ and $\phi = \mathbf{1}_n$, $\sigma^{\eta} = n^{-1/3}\eta^{2/3}$, if $\lim_{\eta \to \infty} \tau_{\max}^{\eta} / \eta^{1/3} = 0$, we have

$$\limsup_{\eta \to \infty} \frac{L^{\eta}}{\eta^{1/3}} \le \left(\sum_{i \in N} \left(\frac{\mu_{i}^{*} G_{i}^{"}(\mu_{i}^{*})}{2} + G_{i}^{'}(\mu_{i}^{*}) \right) - \sum_{j \in N} \left(\frac{\lambda_{j}^{*} F_{j}^{"}(\lambda_{j}^{*})}{2} + F_{j}^{'}(\lambda_{j}^{*}) \right) + \frac{\max_{i \in N, j \in M} \{ s_{i}^{(s)}, s_{j}^{(c)} \}}{2} (\langle \mathbf{1}_{n}, \boldsymbol{\lambda}^{*} \rangle + \langle \mathbf{1}_{n}, \boldsymbol{\mu}^{*} \rangle) \right) n^{-2/3}$$

$$= O(n^{1/3})$$

4.3. Lower Bound

In this section, we prove a lower bound on the profit loss for the large-system regime. We show that under mild assumptions, any pricing and matching policy has a profit loss of $\Omega((n\eta)^{1/3})$. In light of this lower-bound result and the upper bound from Theorem 5, the two-price max-weight policy achieves the optimal rate for the large-system regime.

To prove the lower bound, we consider the following problem instance. We assume that the numbers of customer and server types are equal (i.e., n = m) and that the bipartite matching graph is a complete graph. Intuitively, the complete graph presents a best-case scenario, as it provides with the maximum flexibility to match any customer-server pairs. To define the supply and demand rate, we assume that there exist functions $F, G : \mathbb{R}_+ \to \mathbb{R}_+$ such that $F_j = F$ for all $j \in M$ and $G_i = G$ for all $i \in N$. It is easily verified that this problem instance satisfies the complete resource pooling condition (Condition 2). We also assume that the unit waiting cost is $\mathbf{s} = \mathbf{1}_{2n}$. By symmetry, we can conclude that the fluid solution $(\boldsymbol{\lambda}^{\star}, \boldsymbol{\mu}^{\star}) = (\lambda^{\star} \mathbf{1}_{n}, \mu^{\star} \mathbf{1}_{n})$ for some $\lambda^* = \mu^* > 0$. We show the following properties of the optimal pricing and matching policy for this instance.

Proposition 7. Consider a two-sided queueing system defined by a complete bipartite graph $G(M \cup N, M \times N)$. Assume that n = m, $F_j = F$ for all $j \in M$, $G_i = G$ for all $i \in N$,

and $\mathbf{s} = \mathbf{1}_{2n}$. Then, there exists an optimal policy $(\boldsymbol{\lambda}^{\star}(\cdot), \boldsymbol{\mu}^{\star}(\cdot), \mathbf{x}^{\star}(\cdot))$ that satisfies the following:

a. $\lambda_{j_1}^{\star}(\mathbf{q}) = \lambda_{j_2}^{\star}(\mathbf{q}) \text{ and } \mu_{i_1}^{\star}(\mathbf{q}) = \mu_{i_2}^{\star}(\mathbf{q}) \text{ for all } j_1, j_2 \in M,$ $i_1, i_2 \in N \text{ and } \mathbf{q} \in \mathbb{Z}_+^{2n};$

b. $(\boldsymbol{\lambda}^{\star}(\mathbf{q}_1), \boldsymbol{\mu}^{\star}(\mathbf{q}_1)) = (\boldsymbol{\lambda}^{\star}(\mathbf{q}_2), \boldsymbol{\mu}^{\star}(\mathbf{q}_2))$ if $\langle \mathbf{1}_n, \mathbf{q}_1^{(c)} \rangle = \langle \mathbf{1}_n, \mathbf{q}_2^{(c)} \rangle$ and $\langle \mathbf{1}_n, \mathbf{q}_1^{(s)} \rangle = \langle \mathbf{1}_n, \mathbf{q}_2^{(s)} \rangle$; and

c. $\langle \mathbf{1}_n, \mathbf{x}^*(\mathbf{q}) \rangle = 2\min\{\overline{\langle \mathbf{1}_n, \mathbf{q}^{(c)} \rangle}, \langle \mathbf{1}_n, \mathbf{q}^{(s)} \rangle\}$ for all $\mathbf{q} \in \mathbb{Z}_+^{2n}$.

Part (a) implies that the optimal prices are equal for all the customer queues and server queues, respectively. Part (b) implies that the arrival rates depend on the state of the system **q** only through the total number of customers and total number of servers in the system. Both of these conclusions are intuitive as the problem instance is defined symmetrically for all types. Lastly, part (c) implies that the optimal policy will always match the maximum possible number of customerserver pairs. In particular, there is no incentive to hold a compatible customer-server pair because of the complete graph structure.

Motivated by the proposition, we restrict ourselves to a family of *symmetric pricing policies* defined as follows:

$$\lambda_{j}^{\eta}(\mathbf{q}) = \eta \lambda^{*} + f\left(\frac{\langle \mathbf{1}_{n}, \mathbf{q}^{(c)} \rangle - \langle \mathbf{1}_{n}, \mathbf{q}^{(s)} \rangle}{\eta^{\alpha}}\right) \eta^{\beta} \quad \forall j \in M,$$
(13)

$$\mu_{i}^{\eta}(\mathbf{q}) = \eta \mu^{*} + g \left(\frac{\langle \mathbf{1}_{n}, \mathbf{q}^{(c)} \rangle - \langle \mathbf{1}_{n}, \mathbf{q}^{(s)} \rangle}{\eta^{\alpha}} \right) \eta^{\beta} \quad \forall i \in N.$$
(14)

The policy family is based on the one considered in Section 3.4. However, we make two changes on Equations (10) and (11). We assume $f_j = f$ ($\forall j \in M$) and $g_i = g$ ($\forall i \in N$) because the prices are symmetric for all customer and server types by part (a); the policy depends on the state \mathbf{q} through the difference between total customers and total servers, according to parts (b) and (c). We impose some technical conditions on f and g.

Condition 4.

a. There exists a constant Γ , which may depend on n but not on η , such that $|f(z)| \le \Gamma$ and $|g(z)| \le \Gamma$ for all $z \in \mathbb{R}$.

b. $0 < \alpha + \beta < 1$.

c. There exist constants $\kappa, \delta > 0$, such that if $z > \kappa$, then either $f(z) < -\delta\Gamma$ or $g(z) > \delta\Gamma$. In addition, if $z < -\kappa$, then either $f(z) > \delta\Gamma$ or $g(z) < -\delta\Gamma$.

Condition 4(a) is analogous to Condition 1(a), implying that f and g are bounded. Condition 4(b) is similar to Condition 1(b). Lastly, Condition 4(c) is adapted from Condition 1(c) using the properties from Proposition 7. Now, we present the result on the lower bound.

Theorem 6. Consider a family of two-sided queues parametrized by η and n in the large-system regime. Under any symmetric pricing policy satisfying Condition 4 and any matching policy, we have

$$\liminf_{\eta \to \infty} \frac{L^{\eta}}{\eta^{1/3}} = \Omega(n^{1/3}).$$

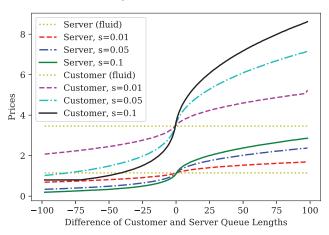
Remark 3 (Intuitive Explanation of $n^{1/3}$). We prove the lower bound by considering the complete graph instance with symmetric demand and supply functions. By Proposition 7, the optimal policy will match the maximum possible number of customer-server pairs, so only the excess customers/servers are waiting in the system. Thus, the system mimics a single-link twosided queue (m = n = 1) with customer and server arrival rates given by $\langle 1, \lambda \rangle$ and $\langle 1, \mu \rangle$, respectively. Under the large-system regime where the number of types is scaled by *n* and the arrival rate of each type is scaled by η , the total arrival rates for customers and servers are scaled by $n\eta$. The profit loss of any pricing and matching policy is lower bounded by the cube root of the total arrival rate (see Remark 2), so the profit loss in the large-system regime is of order $\Omega((n\eta)^{1/3})$. Meanwhile, because the max-weight matching policy leads to state space collapse under the CRP condition, the two-sided queueing system behaves like a singlelink system where all customer and server types are pooled together, so the two-price max-weight policy is able to achieve a tight $O((n\eta)^{1/3})$ profit loss rate.

5. Numerical Experiments 5.1. Single-Link Systems

Our first experiment analyzes a single-link system with one server type and one customer type. In this case, the system state of the MDP is represented by a single variable: namely, the difference between the customer queue length and the server queue length (a detailed discussion of this system is included in Online Appendix A.2). We will solve the optimal policy of the MDP and compare it with the fluid pricing policy and the two-price policy.

We assume a supply curve given by $p_1 = \lambda^{0.5}$ and a demand curve given by $p_2 = 4\mu^{-0.5}$. With these supply and demand curves, the optimal profit of the fluid model is 3.08 when $\lambda = \mu = 4/3$, $p_1 = 1.15$, and $p_2 = 3.46$. We then calculate the optimal pricing policy of the long-run average cost MDP using relative value iteration. Figure 3 shows the optimal pricing policy under three different values of the penalty coefficient (s), as well as the optimal price of the fluid model. The result shows that the optimal customer price is always above the server price, and both prices are increasing with the queue-length difference. Intuitively, if the system has more customers, the customer price should be

Figure 3. (Color online) Optimal Pricing Policies Under Different Values of Penalty Coefficients



increased to reduce the customer arrival rate, and the server price should be increased to increase the server arrival rate. This observation verifies Proposition EC.1 in Online Appendix A.2. As *s* increases, more weight is given to the waiting cost (or equivalently, customers and servers become more sensitive to delays), so the price increases more steeply as the numbers of customers and servers waiting in the system increase. Figure 4 shows the stationary distribution and the mean of queue length for different values of the penalty coefficient (*s*). As expected, when *s* increases, the queue length is more concentrated around zero.

Furthermore, we simulate the profit loss under the fluid pricing policy and two-price policy and compare it with the theoretical result presented before and also with the exact solution obtained by solving the MDP. The result is presented in Figure 5. The profit loss under the fluid pricing policy has an order of $\sqrt{\eta}$, and that under the two-price policy has an order of $\eta^{1/3}$, verifying Theorems 1 and 2. Also, observe that the profit loss under the two-price policy is not much

Figure 4. (Color online) Stationary Distribution of Queue Length Under Different Penalty Coefficients

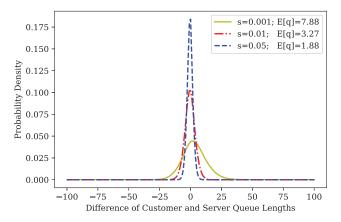
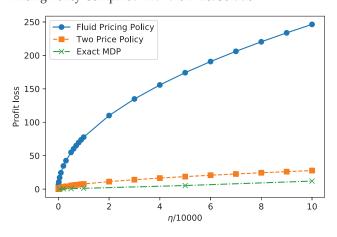


Figure 5. (Color online) Performance of Two-Price and Fluid Pricing Policy Compared with the Exact Solution



different from that of the optimal profit loss, demonstrating the effectiveness of a two-price policy.

5.2. Systems with Multiple Types

Next, we analyze the general two-sided queues with multiple customer and server types. We examine the profit losses under the following four different algorithms.

- 1. FP + MW represents the fluid pricing (Equation (7)) and max-weight matching (Equation (5)) policy.
- 2. FP + Rand represents the fluid pricing (Equation (7)) and randomized matching (Algorithm 2) policy.
- 3. TP + MW represents the two-price policy (Equation (9)) with max-weight matching (Equation (5)).
- 4. TP + Rand represents the two-price policy (Equation (9)) with randomized matching (Algorithm 2).

In this numerical experiment, we first consider a setting where the number of servers and the number of customers are equal (m = n) and where CRP condition (Condition 2) is satisfied. We assume the compatibility graph is given by

$$E = \{(i,j) \in [n] \times [n] : j \in \{i+k\} \cup \{(i+k-n)^+\},\$$

$$k = 0,1,2,3\}.$$

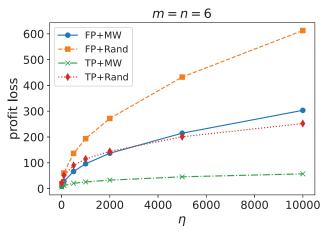
The demand and supply curves are given by

$$F_j(\lambda_j) = 2 - \lambda_j/2, \ \forall j \in [m],$$
 and $G_i(\mu_i) = \mu_i/2 \ \forall j \in [n],$

respectively. We assume the unit holding cost is s=1. The parameter of the fluid pricing policy is set to $q_{\max}^{\eta}=2\sqrt{\eta/n}$. The parameters of the two-price policy are chosen to be $\tau_{\max}^{\eta}=0$ and $\sigma^{\eta}=\eta^{2/3}n^{-1/3}$.

We report the profit loss for $\eta \in \{10, 100, 500, 1, 000, 2, 000, 5, 000, 10, 000\}$ when m = n = 6 in Figure 6. We find that when η is larger, the profit loss of TP + MW grows the slowest, followed by the profit loss of TP + Rand, FP + MW, and FP + Rand. This result confirms

Figure 6. (Color online) Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different η When m = n = 6



the advantage of the two-price policy over the fluid pricing policy, as well as the advantage of the maxweight matching policy over the randomized matching policy. Figure 7 shows the same plot in logarithmic scale. Note that the slope of the log-log plot in Figure 7 can be interpreted as the order of profit loss with respect to η . The fitted slopes of FP + MW and TP + MW are 0.51 and 0.33, respectively. This is consistent with Theorems 1 and 2, which state that FP + MW and TP + MW have the orders of profit loss with respect to η of $O(\eta^{1/2})$ and $O(\eta^{1/3})$, respectively. Figure 7 shows that FP + MW and FP + Rand yield the same order of profit loss with respect to η of approximately 1/2. Moreover, the two-price policy combined with either max-weight matching or randomized matching yields the same order of profit loss with respect to η of approximately 1/3. That is, choosing max-weight or randomized matching does not affect order of profit loss with respect to η .

Figure 7. (Color online) Log-Log Plot of Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different η When m = n = 6

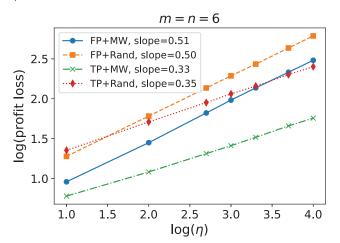
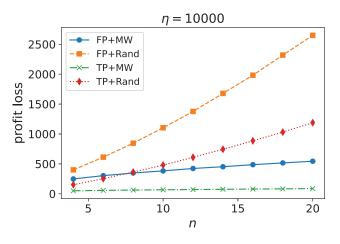
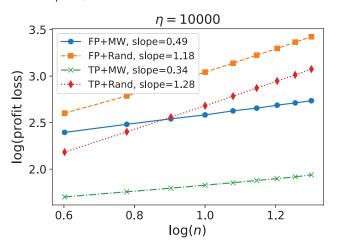


Figure 8. (Color online) Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different n When $\eta = 10.000$



Our next experiment investigates how the profit loss changes with the number of customer and server types. We first consider a setting when the types of two sides are balanced (m = n). Figure 8 shows the profit loss for $n \in \{4, 6, 8, ..., 20\}$ when $\eta = 10,000$ (a large η is chosen so that the asymptotic trend becomes clear). Figure 9 shows the same plot in logarithmic scale. It can be observed that the profit losses when a pricing policy is combined with the randomized matching policy grow faster than those when a pricing policy is combined with the max-weight matching policy as *n* increases. In other words, the max-weight matching policy performs better than the randomized matching policy. Figure 9 suggests that the orders of profit loss with respect to n of FP + MW and FP + Rand are 0.49 and 1.18, respectively, which are close to those predicted by Theorem 4. Moreover, the orders of profit loss with respect to n of TP + MW and TP +

Figure 9. (Color online) Log-Log Plot of Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different n When $\eta = 10,000$



Rand are 0.34 and 1.28, respectively, which are close to those predicted by Theorem 5. Clearly, in both cases, the max-weight algorithm performs much better than the randomized matching algorithm for large n.

We also consider the setting where the number of server queues and the number of customer queues are not equal. Specifically, we assume that the number of server queues is twice as many as the number of customer queues (i.e., m = 2n). The compatibility graph is given by

$$E = \{(i,j) \in [2n] \times [n] : j \in \{i+k\} \cup \{(i+k-n)^+\}, k = 0,1\}.$$

The demand and supply curves are assumed to be

$$F_j(\lambda_j) = 6 - \lambda_j, \quad \forall j \in [m]$$
 and $G_i(\mu_i) = \mu_i, \quad \forall i \in [n],$

respectively. The parameters of pricing policy are similar to the previous case when m = n.

We report the profit loss for $\eta \in \{10, 100, 500,$ 1,000,2,000,5,000,10,000} when m = 8 and n = 4 in Figure 10. The result shows that when η is larger, the profit loss when using the two-price policy grows significantly slower, compared with when using the fluid pricing policy. Moreover, we can observe that in this case, the benefit of the max-weight matching policy over the randomized matching policy when combined with any pricing policy is negligible. The same plot in logarithmic scale (shown in Figure 11) shows that the fitted orders of profit loss with respect to η of FP + MW and FP + Rand are 0.49 and 0.48, respectively, and that those of TP + MW and TP + Rand are 0.33 and 0.32, respectively. This observation confirms the results from Theorems 1 and 4 as well as the results from Theorems 2 and 5, which state that the orders of

Figure 10. (Color online) Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different η When m=8 and n=4

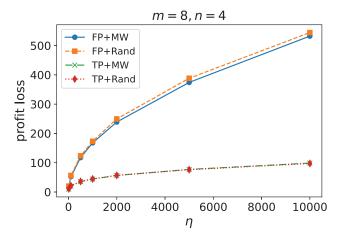
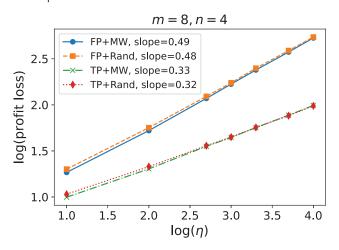


Figure 11. (Color online) Log-Log Plot of Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different η When m = 8 and n = 4



profit loss with respect to η of FP + MW and FP + Rand are 1/2 and that those of TP + MW and TP + Rand are 1/3, respectively. Figure 12 shows the profit loss for $n \in \{4,6,8,\ldots,20\}$ and m=2n when $\eta=10,000$. Figure 13 shows the same plot in logarithmic scale. These figures show the superiority of the max-weight policy over the randomized matching policy discussed in Section 4.2.

6. Conclusion

In this paper, we present a model of dynamic pricing and matching for two-sided queueing systems. The system is formulated as a Markov decision process, and a fluid approximation model is considered. We presented a fluid pricing and max-weight matching policy and showed that it achieves $O(\sqrt{\eta})$ optimality rate. Furthermore, we proposed a dynamic pricing and max-weight policy, which achieves $O(\eta^{1/3})$ optimality rate. We also show that this scaling of $O(\eta^{1/3})$

Figure 12. (Color online) Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different n When $\eta = 10,000$

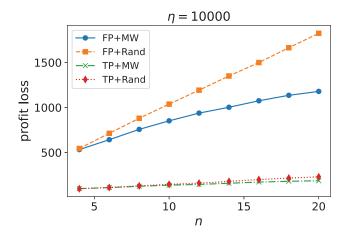
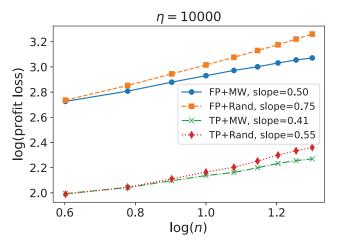


Figure 13. (Color online) Log-Log Plot of Profit Losses Under FP + MW, FP + Rand, TP + MW, and TP + Rand for Different n When $\eta = 10,000$



matches the lower bound for a broad family of policies. We also demonstrate the advantage of maxweight matching over randomized matching. Under the complete resource pooling condition, we show that max-weight matching achieves $O(\sqrt{n})$ and $O(n^{1/3})$ optimality rates for static and two-price policies, respectively, where n is the number of customer and server types. In comparison, the randomized matching policy may have an $\Omega(n)$ optimality rate.

References

Adan I, Weiss G (2012) Exact FCFS matching rates for two infinite multitype sequences. *Oper. Res.* 60(2):475–489.

Akbarpour M, Li S, Gharan SO (2020) Thickness and information in dynamic matching markets. *J. Political Econom.* 128(3):783–815.

Anderson R, Ashlagi I, Gamarnik D, Kanoria Y (2017) Efficient dynamic barter exchange. *Oper. Res.* 65(6):1446–1459.

Armony M, Shimkin N, Whitt W (2009) The impact of delay announcements in many-server queues with abandonment. *Oper. Res.* 57(1):66–81.

Banerjee S, Freund D, Lykouris T (2017) Pricing and optimization in shared vehicle systems: An approximation framework. Daskalakis C (General Chair), ed. *Proc.* 2017 ACM Conf. Econom. Comput. (Association for Computing Machinery, New York), 517. https://dl.acm.org/doi/abs/10.1145/3033274.3085099.

Banerjee S, Johari R, Riquelme C (2016) Dynamic pricing in ride-sharing platforms. *ACM SIGecom Exchanges* 15(1):65–70.

Banerjee S, Kanoria Y, Qian P (2018) State dependent control of closed queueing networks. *Performance Evaluation Rev.* 46(1):2–4.

Besbes O, Castro F, Lobel I (2021) Surge pricing and its spatial supply response. *Management Sci.* 67(3):1350–1367.

Caldentey R, Kaplan EH, Weiss G (2009) FCFS infinite bipartite matching of servers and customers. Adv. Appl. Probab. 41(3): 695–730.

Chen H, Frank MZ (2001) State dependent pricing with a queue. *IIE Trans.* 33(10):847–860.

Chen Y, Hu M (2020) Pricing and matching with forward-looking buyers and sellers. *Manufacturing Service Oper. Management* 22(4):717–734.

Eryilmaz A, Srikant R (2012) Asymptotically tight steady-state queue length bounds implied by drift conditions. *Queueing Systems* 72(3-4):311–359.

- Gallego G, Van Ryzin G (1994) Optimal dynamic pricing of inventories with stochastic demand over finite horizons. *Management Sci.* 40(8):999–1020.
- Gamarnik D, Zeevi A (2006) Validity of heavy traffic steady-state approximations in generalized Jackson networks. Ann. Appl. Probab. 16(1):56–90.
- Gurvich I, Ward A (2014) On the dynamic control of matching queues. *Stochastic Systems* 4(2):479–523.
- Gurvich I, Whitt W (2009) Scheduling flexible servers with convex delay costs in many-server service systems. Manufacturing Service Oper. Management 11(2):237–253.
- Harrison JM (2013) Brownian Models of Performance and Control (Cambridge University Press, Cambridge, UK).
- Hu M, Zhou Y (2021) Dynamic type matching. *Manufacturing Service Oper. Management*, ePub ahead of print March 18, https://pubsonline.informs.org/doi/abs/10.1287/msom.2020.0952.
- Hu B, Hu M, Zhu H (2021) Surge pricing and two-sided temporal responses in ride-hailing. *Manufacturing Service Oper. Manage*ment, ePub ahead of print February 3, https://doi.org/10.1287/ msom.2020.0960.
- Kanoria Y, Qian P (2019) Near optimal control of a ride-hailing platform via mirror backpressure. Preprint, submitted March 7, https://arxiv.org/abs/1903.02764v1.
- Kim J, Randhawa RS (2017) The value of dynamic pricing in large queueing systems. *Oper. Res.* 66(2):409–425.
- Lange DH, Maguluri ST (2019) Heavy-traffic analysis of the generalized switch under multidimensional state space collapse. *Perfor*mance Evaluation Rev. 47(2):36–38.
- Low DW (1974a) Optimal dynamic pricing policies for an M/M/s queue. *Oper. Res.* 22(3):545–561.
- Low DW (1974b) Optimal pricing for an unbounded queue. *IBM J. Res. Development* 18(4):290–302.
- Maguluri ST, Srikant R (2015) Queue length behavior in a switch under the maxweight algorithm. Preprint, submitted June 10, https://arxiv.org/abs/1503.05872.
- Maguluri ST, Srikant R (2016) Heavy traffic queue length behavior in a switch under the MaxWeight algorithm. *Stochastic Systems* 6(1):211–250
- Nguyen LM, Stolyar AL (2018) A queueing system with on-demand servers: Local stability of fluid limits. *Queueing Systems* 89(3-4): 243–268.
- Özkan E, Ward AR (2020) Dynamic matching for real-time ride sharing. *Stochastic Systems* 10(1):29–70.
- Paschalidis IC, Tsitsiklis JN (2000) Congestion-dependent pricing of network services. IEEE/ACM Trans. Networking 8(2):171–184.
- Puterman ML (1994) Markov Decision Processes: Discrete Stochastic Dynamic Programming (John Wiley & Sons, Inc., New York).
- Roth AE, Sönmez T, Ünver MU (2007) Efficient kidney exchange: Coincidence of wants in markets with compatibility-based preferences. Amer. Econom. Rev. 97(3):828–851.
- Shi C, Wei Y, Zhong Y (2019) Process flexibility for multiperiod production systems. *Oper. Res.* 67(5):1300–1320.
- Sivaraman V, Venkatakrishnan SB, Ruan K, Negi P, Yang L, Mittal R, Fanti G, Alizadeh M (2020) High throughput cryptocurrency routing in payment channel networks. Bhagwan

- R (Chair), Porter G, eds. 17th USENIX Sympos. Networked Systems Design Implementation (NSDI 20), (USENIX Association, Berkeley, CA), 777–796. https://www.usenix.org/conference/nsdi20/presentation/sivaraman.
- Srikant R, Ying L (2014) Communication Networks: An Optimization, Control and Stochastic Networks Perspective (Cambridge University Press, New York).
- Stolyar AL (2004) MaxWeight scheduling in a generalized switch: State space collapse and workload minimization in heavy traffic. Ann. Appl. Probab. 14(1):1–53.
- Tan B, Srikant R (2012) Online advertisement, optimization and stochastic networks. *IEEE Trans. Automatic Control* 57(11): 2854–2868.
- Tassiulas L, Ephremides A (1992) Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks. *IEEE Trans. Automatic Control* 37(12):1936–1948.
- Williams RJ (1998) Diffusion approximations for open multiclass queueing networks: Sufficient conditions involving state space collapse. Queueing Systems 30(1-2):27–88.
- Yan C, Zhu H, Korolko N, Woodard D (2020) Dynamic pricing and matching in ride-hailing platforms. *Naval Res. Logist.* 67(8): 705–724.
- Zohar E, Mandelbaum A, Shimkin N (2002) Adaptive behavior of impatient customers in tele-queues: Theory and empirical support. Management Sci. 48(4):566–583.

Sushil Mahavir Varma is currently pursuing a PhD in operations research, supervised by Prof. Siva Theja Maguluri in the H. Milton Stewart School of Industrial and Systems Engineering at Georgia Tech. He is currently working on problems spanning queueing theory, game theory, and its applications in ride hailing, blockchain, wireless networks, etc.

Pornpawee Bumpensanti got a PhD degree in operations research from Georgia Institute of Technology. Her research focuses on pricing and revenue management problems.

Siva Theja Maguluri is Fouts Family Early Career Professor and an assistant professor in the School of Industrial and Systems Engineering at Georgia Tech. His research interests span the areas of networks, control, optimization, algorithms, applied probability, and reinforcement learning. He is a recipient of NSF CAREER award, "Best Publication in Applied Probability", and "Student Recognition of Excellence in Teaching" Award.

He Wang is an assistant professor and the Colonel John B. Day Early Career Professor in the H. Milton Stewart School of Industrial and Systems Engineering at Georgia Tech. His research focuses on the interface between machine learning and operations management, where he develops data-driven methods for applications in dynamic pricing, supply chain, transportation, and marketplace design.