

Singapore Management University

Institutional Knowledge at Singapore Management University

Research Collection Lee Kong Chian School Of
Business

Lee Kong Chian School of Business

11-2023

Maximizing the benefits of an on-demand workforce: Fill rate-based allocation and coordination mechanisms

Tao LU

University of Connecticut - Storrs

Zhichao ZHENG

Singapore Management University, DANIELZHENG@smu.edu.sg

Yuanguang ZHONG

South China University of Technology

Follow this and additional works at: https://ink.library.smu.edu.sg/lkcsb_research



Part of the [Operations and Supply Chain Management Commons](#), and the [Strategic Management Policy Commons](#)

Citation

LU, Tao; ZHENG, Zhichao; and ZHONG, Yuanguang. Maximizing the benefits of an on-demand workforce: Fill rate-based allocation and coordination mechanisms. (2023). *Manufacturing and Service Operations Management*. 25, (6), 2216-2232.

Available at: https://ink.library.smu.edu.sg/lkcsb_research/7263

This Journal Article is brought to you for free and open access by the Lee Kong Chian School of Business at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection Lee Kong Chian School Of Business by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email cherylds@smu.edu.sg.

Maximizing the Benefits of an On-Demand Workforce: Fill Rate-Based Allocation and Coordination Mechanisms

Tao Lu¹, Zhichao Zheng², and Yuanguang Zhong³

¹School of Business, University of Connecticut, Storrs, Connecticut, USA

²Lee Kong Chian School of Business, Singapore Management University, Singapore

³School of Business Administration, South China University of Technology, Guangzhou, China

Abstract

Problem definition: With the rapid growth of the gig economy, on-demand staffing platforms have emerged to help companies manage their temporary workforce. This emerging business-to-business context motivates us to study a new form of supply chain coordination problem. We consider a staffing platform managing an on-demand workforce to serve multiple firms facing stochastic labor demand. Before demand realization, each individual firm can hire permanent employees, whereas the platform determines a compensation rate for potential on-demand workers. After knowing the realized demand, firms in need can request on-demand workers from the platform, and then the platform operator allocates the available on-demand workforce among the firms. We explore how to maximize and distribute the benefits of an on-demand workforce through coordinating self-interested parties in the staffing system. **Methodology/results:** We combine game theory and online optimization techniques to address the challenges in incentivizing and coordinating the online workforce. We propose a novel and easily implementable fill rate-based allocation and coordination mechanism that enables the on-demand workforce to be shared optimally when individual firms and the platform operator make decisions in their own interest. We also show that the proposed mechanism can be adapted to the cases when contract terms need to be identical to all firms and when actual demand is unverifiable. **Managerial implications:** The proposed contract mechanism is in line with the performance-based contracting commonly used in on-demand staffing services. Our results suggest that under an appropriately designed performance-based mechanism, individual firms and the platform operator can share the maximum benefits of on-demand staffing.

1 Introduction

The Fourth Industrial Revolution is fundamentally transforming our lives and work through innovative technologies. This shift is also evident in the staffing industry, where online platforms are revolutionizing workforce recruitment and management. The 2020 Gig Economy and Talent Platform Landscape report by Staffing Industry Analysts (SIA, 2020) reveals a 42% revenue increase in online staffing firms in 2019. Besides the growth in conventional business-to-consumer

(B2C) platforms like Uber and Lyft, the business-to-business (B2B) segment is steadily emerging, with 267 B2B-focused talent platforms generating \$9.4 billion in 2019. In 2020, the gross service volume at B2B-focused talent platforms grew by 25% globally, according to an updated report by SIA (2021b).

The B2B online staffing platforms deploy on-demand temporary workers to businesses to fill their hourly or daily labor needs. For example, Jitjatjo provides on-demand workers to various businesses, e.g., catering servers and dishwashers to restaurants, general cleaners, and janitors to commercial facilities in the US and Australia; Instaff offers retail staff (e.g., cashiers and stock replenishers), among many other roles, to firms in Germany; and Weploy provides customer service and general administration roles to businesses in Australia.¹ The ongoing global pandemic of coronavirus disease 2019 (COVID-19) has significantly accelerated the growth of the B2B online staffing platform. Specifically, the on-demand platforms for nurse staffing have been “a lifesaver” to catch the increasing and unpredictable demand for healthcare workers (SIA, 2021a). For example, Jitjatjo started to provide dishwashers and janitors to healthcare facilities, and it further expanded its range of services to supply patient transporters and disinfection technicians, etc.² Besides sourcing the workforce supplies at the tactic and operational levels, companies have started integrating the on-demand workforce into their organizational structures at strategic levels to increase flexibility and allow unified workforce management.

Unlike traditional job-advertisement platforms, these online staffing platforms manage employment online and price via marketplace mechanics (InStaff, 2016, p. 7).³ They prescreen potential workers, allocate available workers to temporary positions posted by customers (i.e., employers) on an on-demand basis, and process payroll and insurance for their customers.

Compared to the conventional staffing channel, in which employers must hire a sufficiently large regular workforce to meet demand in peak periods, the availability of an on-demand workforce allows employers to meet peak demand with a relatively low level of regular employees, thus reducing potential overstaffing costs during off-peak periods. However, the risk of not getting sufficient manpower from an uncertain on-demand workforce poses a significant challenge to employers in terms of how to ration staffing levels between these two channels. Gurvich et al. (2019) investigated the setting of a single employer who must motivate the on-demand workforce via the compensation rate and found that the employer’s staffing level and profit would decrease

¹<https://www.jitjatjo.com>, <https://en.instaff.jobs/sales-staff-germany>, <https://www.weployapp.com/>.

²<https://www.jitjatjo.com/ondemand/healthcare-staffing>.

³In practice, most on-demand staffing platforms determine wage rates for temporary workers and charge their customers—i.e., employers—a service fee. For example, InStaff and Weploy have hourly wage rates posted online; JitJatJo offers localized pricing to ensure the best deal for businesses. See <https://en.instaff.jobs/costs-and-benefits>, <https://au.jitjatjo.com/ondemand/staffing>, <https://www.weployapp.com/weployer/pricing>.

compared to the conventional mode, by which the employer can order any number of agents to work, even at the same compensation rate. The problem would be even more challenging for the aforementioned platforms, which serve *multiple* employers and fill their demand with a common pool of on-demand employees, in the sense that such platforms must design proper allocation and contract mechanisms to coordinate and fulfill the staffing needs of multiple self-interested employers. In this paper, we aim to address the following questions:

- (1) *When is it beneficial for employers to use an on-demand workforce, and to what extent?*
- (2) *Given the fact that the platform operator and employers are independent businesses and make decisions in their own interests, can we maximize the benefits of an on-demand workforce through carefully designed allocation and contract mechanisms? If so, can the mechanism lead to a win-win outcome in which the platform and employers are all better off with on-demand staffing?*
- (3) *An additional challenge arises if employers' actual demands are not verifiable. Can we ensure each employer truthfully report their demand under some coordination mechanism?*

To this end, we study a system that consists of multiple firms as employers and an on-demand staffing platform with a pool of self-scheduling workers. Firms can either recruit permanent employees from a traditional staffing channel or hire on-demand workers via the platform to satisfy their labor demand. On-demand workers have heterogeneous preferences for work, which are uncertain to the platform operator. Therefore, the platform operator must encourage on-demand workers to work by offering an appropriate compensation rate. On the other hand, the platform operator also needs to design appropriate contracts to engage the firms to use the platform's service. The design of the contract is further complicated by how the platform operator allocates on-demand workers to satisfy the staffing needs of the firms. In such a system, employers, as independent business owners, would choose staffing strategies to suit their own interests; at the same time, the platform operator aims to maximize its profit. The incentives of different parties are affected by both the service contracts and the allocation policy. The problem becomes more challenging when firms' actual demands are unverifiable. Despite the increasingly prevalent collaborations between on-demand platforms and employers, to the best of our knowledge, there appears to be no systematic study of the incentive issues for this one-to-many staffing system. In this paper, we focus on developing the optimal mechanism—including contracts and workforce allocation policy—that induces all self-interested players to make system-wide optimal decisions, thereby maximizing the benefits of an on-demand workforce.

To understand the first-best solution, we characterize system-wide optimal staffing strategies

whereby the benefits of an on-demand workforce are maximized from the perspective of the overall system. This also tells us when it is beneficial to use an on-demand workforce compared to a traditional workforce. For the general problem, we first explicitly address workforce allocation among multiple employers, because it determines the actual service level, measured by fill rate⁴—proportion of demand filled by supply—delivered to each customer, which in turn influences each employer’s staffing strategy, i.e., to what extent they would hire on-demand workers. Using techniques from online optimization, we characterize the optimal target fill rates that should be delivered to each employer, and show that these desired fill rates can be implemented under a proper allocation policy. Next, we propose a novel contract mechanism contingent on actual job fulfillment relative to the designed fill rate, and demonstrate that the proposed mechanism induces the system-wide optimal outcome. Moreover, we discuss how the proposed mechanism can be adapted to the cases when contract terms need to be identical to all firms and when actual demand is unverifiable.

Our study thus makes the following contributions: (1) To our knowledge, this is the first paper to study service contracts for on-demand staffing platforms that serve multiple employers. We characterize the conditions under which it is optimal for a one-to-many staffing system to use (or partially use) an on-demand workforce. (2) We propose a novel fill rate-based allocation and contract mechanism to coordinate self-interested parties and induce system-wide optimal solutions. Our results have profound managerial implications: it offers guidance for the use of fill rate—one of the most important performance metrics widely used by staffing companies (Taylor, 2017)—as a protocol for workforce allocation and incentive alignment. Broadly speaking, our work contributes to the capacity pooling literature in that we demonstrate how the risk-pooling benefit can be maximized and redistributed to all self-interested parties, although the literature has often observed that risk pooling may not benefit everyone in various decentralized supply chain settings (e.g., Anupindi and Bassok, 1999; Dong and Rudi, 2004). (3) Additionally, the proposed mechanism after some slight modification provides a novel approach to ensure employers truthfully report their demand, in case actual demands are unverifiable.

2 Related Literature

The literature on operations management in the sharing economy is rapidly growing. For instance, Gurvich et al. (2019) consider a service platform that hires self-scheduling agents to satisfy customers’ demand. They model the service provider as a newsvendor who cannot dic-

⁴Staffing companies view fill rate as an essential metric for understanding their firm’s efficiency. If the fill rate begins to decline, this is taken as a warning sign (Taylor, 2017).

tate the number of workers, but must offer a compensation rate such that agents will decide whether to work based on their individual availability. Cachon et al. (2017) consider a similar newsvendor model, in which the platform can dynamically adjust the price to customers and wages to workers. Taylor (2018) considers the situation in which customers' utility (and thus the demand rate) is affected by the congestion level of the service system, and examines the impact of such system dynamics on the platform operator's optimal price and wage. Chen and Hu (2019) further consider a platform to match multiple supply types to multiple demand types over a planning horizon, in which the optimal match policy is derived based on the priorities of demand-supply pairs. Benjaafar et al. (2020), Benjaafar et al. (2021) and Lin et al. (2022) explore the welfare implications under the on-demand service context. We refer interested readers to Benjaafar and Hu (2019) and references therein. In this paper, we adopt a newsvendor-type platform with self-scheduling workers, as in Gurvich et al. (2019) and Cachon et al. (2017). However, unlike their ride-hailing setting, we model the platform's customers as independent, cost-minimizing businesses motivated by B2B staffing practices. We identify optimal staffing strategies and examine how to coordinate self-interested firms in such a system.

Our paper is relevant to the literature on service and supply chain contracting. In the literature of supply chain coordination, various contracts such as buy-back (Pasternack, 1985) and revenue sharing (Cachon and Lariviere, 2005) have been studied. We refer interested readers to Cachon's (2003) review of this literature and recent papers such as Chen et al. (2016) and Chen and Lee (2016) and references therein. The supply chain coordination literature primarily focuses on supply chains in which a supplier sells to a newsvendor retailer and the retailer places orders before demand realization. In our problem, it is the platform operator (like a common supplier) who controls the on-demand capacity level, whereas employers (like retailers) determine their own permanent staffing levels and submit job orders after demand is realized. Moreover, we consider a one-to-many staffing system with short-term contracts⁵, whereas existing studies on service contracts examine relatively long-term contracts based on queueing models in different contexts such as call centers (e.g., Hasija et al., 2008; Ren and Zhou, 2008).

The staffing system studied in this paper can be viewed as a two-echelon distribution system consisting of a supplier (analog to the platform) and multiple retailers (analog to the hiring firms) (Section 8.5 in Zipkin, 2000). Inventory (analog to the workforce) can be held at each location, including both the supplier and retailer sites. Existing inventory literature has explored two-echelon inventory systems with various focuses, such as coordination of retailers competing

⁵In health care and humanitarian literature, some studies investigate mechanism design problems for the allocation of randomly arriving resources to recipients (see Zhang et al., 2020, and references therein). Our staffing system, however, is very different from theirs.

in an end market with a deterministic demand function (Chen et al., 2001), performance analysis of echelon inventory policies (Gallego and Zipkin, 1999), and competitive selection of inventory policies by decentralized players (Cachon, 2001). Unlike those papers, we consider a single period model where demand is a general multivariate random variable and the inventory held by the supplier is allocated after demand at each retailer is realized. Thus, closer to our study is the literature on inventory pooling and allocation. It is well-known that using a common pool of inventory yields benefits because of the risk-pooling effect (Eppen, 1979). Some recent papers study how to allocate common capacity from a supplier to fulfill random demand from multiple retailers/customers (e.g., Swaminathan and Srinivasan, 1999; Zhang, 2003; Alptekinoglu et al., 2013; Asadpour et al., 2020). In particular, Zhong et al. (2017) propose a randomized allocation policy to satisfy given fill rate requirements from multiple customers. In their paper, however, the incentive issue arising from self-interested firms is not considered.

Existing papers have observed that due to the incentive conflicts among self-interested parties, inventory pooling systems (in a broad sense) may not operate in a system-wide optimal fashion under various contexts such as lateral transshipments (Dong and Rudi, 2004), contract manufacturing (Ülkü et al., 2007), and pooling purchases (Hu et al., 2013). In particular, Anupindi and Bassok (1999) consider a supply chain consisting of a supplier and two retailers where retailers can jointly hold a centralized inventory. Among many other differences, their setting differs from ours in that retailers place orders before demand realization, rather than requesting allocations in an *on-demand* fashion. Cachon and Lariviere (1999) consider an on-demand allocation of a supplier’s capacity to multiple retailers. However, they focus on the impact of pre-announced allocation rules under price-only contracts and do not address the supply chain coordination problem. More similar to our setting, Netessine and Rudi (2006) allow the supplier and retailers to determine their own inventory levels before demand realization, while, if needed, retailers can receive an allocation of the supplier’s inventory on demand. They focus on wholesale price contracts and characterize the equilibrium of the noncooperative game among supply chain parties. Unlike the above papers, our work considers a staffing system with different model features. More importantly, none of the existing papers has explored coordinating contracts for a decentralized inventory pooling system with on-demand allocation, which is the focus of our study.

3 The Model

We consider a staffing system consisting of a set of n firms (“he”) as potential employers, denoted by $I = \{1, 2, \dots, n\}$, and an on-demand staffing platform (“she”) that hires and provides an on-demand workforce. Throughout the paper, the terms “firm” and “employer” will be used

interchangeably.

We focus on homogeneous workers and a single type of job so that the jobs offered by each firm can be filled by any potential workers in the market. This is a reasonable assumption for staffing services provided by many platforms in practice, which fill nonspecialized job positions such as dishwashers in restaurants and cashiers in retail stores (e.g., Jitjatjo, InStaff). For example, a dishwasher can work for any restaurant, and a cashier can be assigned to work for most retail stores. Firm i ($i \in I$) faces a stochastic demand X_i for the workforce, which can be measured by the number of workers needed for the business.⁶ Let $F_i(x)$ denote the marginal distribution of X_i , and $\bar{F}_i(x) = 1 - F_i(x)$. For technical convenience, we assume that X_i is a continuous and nonnegative random variable with $F_i(0) = 0$, for all $i \in I$. The demands for different firms are not necessarily independent or identical. We assume that the first moment of X_i is finite for all $i \in I$. The distributions F_i are common knowledge to the platform operator and all firms, but individual workers do not need to know this demand information.

Firms can recruit permanent workers from a traditional channel before demand is realized. We denote by Q_i the permanent staffing level of Firm i . Once recruited, permanent workers cannot be laid off. The wage rate for the permanent workforce is exogenously given by c . We assume c to be identical for all firms, since we consider the same type of labor and c represents a standard wage rate in the (permanent) labor market.⁷ Since the staffing level of each firm may not exactly match the demand, an understaffing cost will be incurred at a rate p when the number of workers available for an employer is less than his demand. The value of p can, for instance, represent the overtime pay rate required by labor laws.⁸

Firms may request on-demand workers from a staffing platform to cover personnel shortfalls, denoted by $\hat{\mathbf{X}} = (\hat{X}_1, \hat{X}_2, \dots, \hat{X}_n)$, where $\hat{X}_i = (X_i - Q_i)^+ := \max\{0, X_i - Q_i\}$ for all $i \in I$. For ease of exposition, in the main text we assume that there is a population of infinitesimal self-scheduling workers with size N who decide whether to work via the platform based on their own availability.⁹ In Appendix E, we show that the proposed coordination mechanism remains valid even when the assumption of infinitesimal workers is relaxed. An on-demand worker is willing to work only if his earnings are greater than or equal to his availability threshold. The

⁶For example, if one worker's regular number of hours worked per day is 8, and in a single day Firm i has a task requiring 12 hours of labor, then the firm's demand on that day equals 1.5.

⁷Our results can be readily extended to the case in which firms have heterogeneous staffing cost for hiring permanent workforce.

⁸According to the Fair Labor Standards Act (FLSA) the overtime pay rate for eligible employees is at least one and one-half times their regular wage rate (see <https://www.dol.gov/whd/flsa/>). Therefore, the understaffing cost due to overtime compensation in this case is given by $p = 1.5c$.

⁹By the nature of infinitesimal workers, there are infinitely many workers each of who offer an infinitely small amount of labor supply that adds up to a total supply of N . Such an assumption is quite common in the literature (e.g., Gurvich et al., 2019). In the rest of the paper, more precisely, the notation of one on-demand worker means one unit of on-demand labor supply.

availability threshold τ for each worker is drawn independently from an identical distribution G . We assume that G is continuous and log-concave (Bagnoli and Bergstrom, 2005), and the density function g is strictly positive on a given support $[\underline{\tau}, \bar{\tau}]$. Let $\bar{G}(\tau) = 1 - G(\tau)$. We assume that permanent workers and the on-demand workforce have similar capabilities for the job¹⁰ and the total amount of labor supply from on-demand workers, N , is sufficiently large. To encourage enough workers to work, the platform announces a compensation rate η to potential workers at the beginning,¹¹ and then allocates the available on-demand workforce, denoted by S , to fill the firms' job vacancies $\hat{\mathbf{X}}$. By the assumption of infinitesimal workers, $S = NG(\eta)$ if the compensation rate is set as η . Note that the above model setup implies that in the main paper we will focus on a setting in which workers get paid based on their on-call time, instead of their actual working time. This assumption is in line with Gurvich et al. (2019) and corresponds to the practical situations in which on-demand workers who have agreed to be available are subject to certain restrictions while waiting for assignment.¹² In reality, some platforms may adopt an alternative compensation scheme under which on-demand workers are paid based on their actual working time (see, e.g., Cachon et al., 2017)). In Appendix F, we discuss this alternative setting and show that the proposed coordination mechanism is still applicable.¹³

We denote by $A_i(\hat{\mathbf{X}}, S)$ the number of on-demand workers allocated to Firm i . The allocations $\mathbf{A}(\hat{\mathbf{X}}, S) = (A_1(\hat{\mathbf{X}}, S), A_2(\hat{\mathbf{X}}, S), \dots, A_n(\hat{\mathbf{X}}, S))$ may depend on each firm's job vacancies and the number of available workers. We allow the platform to *randomize* her allocation policy, i.e., $\mathbf{A}(\hat{\mathbf{X}}, S)$ does not need to be a deterministic function of $\hat{\mathbf{X}}$ and S . In general, an allocation policy should satisfy the following restrictions:

$$\sum_{i \in I} A_i(\hat{\mathbf{X}}, S) \leq S, \text{ and } 0 \leq A_i(\hat{\mathbf{X}}, S) \leq \hat{X}_i, \text{ almost surely, } \forall \hat{\mathbf{X}} \in \Omega, \quad (1)$$

where Ω denotes the set of all realizations of $\hat{\mathbf{X}}$, and these conditions must be satisfied almost surely if the allocation policy is not deterministic. That is, the total amount of on-demand workers allocated cannot exceed the available workforce level S , and each firm receives a non-

¹⁰This assumption may depend on the types of jobs, yet it is reasonable for nonspecialized jobs as we consider in the paper.

¹¹On-demand staffing platforms often post the wage rate on their websites. For instance, InStaff and Weploy have hourly wage rates posted online; see <https://en.instaff.jobs/costs-and-benefits>, <https://www.weployapp.com/weployer/pricing>.

¹²For instance, when on-call workers are required to respond quickly to calls or stay within a limited distance from work, they would be entitled to be paid while on call. See a detailed discussion about on-call compensation at https://www.shrm.org/resourcesandtools/tools-and-samples/hr-qa/pages/cms_020208.aspx. Moreover, in some areas it is considered to violate labor laws if on-demand workers will not be paid when they do not get the opportunity to work (Strauss Law Blog, 2015).

¹³When temporary workers are paid based on actual working time, they will decide whether to participate based on the rational expectation of future earnings. This will only complicate the computation of the system-wide optimal solution, but one can still use the form of contracts proposed in our paper to coordinate the system.

negative allocation that is not greater than his request. We define \mathcal{A} as the set of all allocation policies that satisfy (1). To use the on-demand workforce, each firm makes a transfer payment $w_i(\mathbf{A}, \hat{\mathbf{X}})$ to the platform operator, which may generally depend on firms' requests $\hat{\mathbf{X}}$ and the actual allocations $\mathbf{A}(\hat{\mathbf{X}}, S)$.¹⁴ For ease of notation, we will sometimes suppress the dependence of A_i on $(\hat{\mathbf{X}}, S)$ and use $\mathbf{A}(\cdot)$ to emphasize the policy nature of the allocation as a function, and similarly for w_i and $\mathbf{w}(\cdot)$.

In line with the literature on supply chain coordination (see, e.g., Cachon, 2003; Chen et al., 2016), we adopt the Nash equilibrium as our solution concept and assume that all parties are risk-neutral. For any given transfer payment scheme, $\mathbf{w}(\cdot) = (w_1(\cdot), w_2(\cdot), \dots, w_n(\cdot))$, the platform operator chooses the allocation policy $\mathbf{A}(\cdot)$ and the compensation rate η to maximize her expected profit by solving the following problem:

$$\begin{aligned} \max_{\eta, \mathbf{A}(\cdot)} \quad & \Pi(\eta, \mathbf{A}(\cdot) | \mathbf{Q}) = \sum_{i \in I} \mathbb{E} \left[w_i(\mathbf{A}, \hat{\mathbf{X}}) \right] - \eta NG(\eta) \\ \text{s.t.} \quad & \mathbf{A}(\hat{\mathbf{X}}, S) \in \mathcal{A}, S = NG(\eta) \end{aligned} \quad (2)$$

whereas each Firm i determines the permanent staffing level Q_i to minimize his expected staffing cost as follow:

$$\min_{Q_i \geq 0} \quad h_i(Q_i | \mathbf{Q}_{-i}, \mathbf{A}(\cdot), \eta) = \mathbb{E} \left[p \left((X_i - Q_i)^+ - A_i \right)^+ + w_i(\mathbf{A}, (\mathbf{X} - \mathbf{Q})^+) \right] + cQ_i. \quad (3)$$

where \mathbf{Q}_{-i} represents the vector of all the other firms' staffing levels and $((X_i - Q_i)^+ - A_i)^+$ is the final understaffing level of Firm i after using the on-demand workforce. Note that the expectation is taken over both \mathbf{X} and \mathbf{A} if the allocation is randomized. To rule out trivial solutions, we assume $p > c > 0$ such that firms will maintain some permanent workers without any on-demand workforce, and $p > \underline{\tau} > 0$ such that it is profitable to utilize the on-demand workforce without any permanent workers.

The *sequence of events* is summarized as follows. (a) Under a payment scheme $\mathbf{w}(\cdot)$, the platform operator sets and announces her allocation policy $\mathbf{A}(\cdot)$, and compensation rate η ; simultaneously, each firm determines his permanent staffing level Q_i . (b) Knowing the compensation rate, each worker decides whether to work based on an availability threshold τ , drawn from distribution G , and receives compensation η , if available. Meanwhile, each firm's labor demand X_i is realized, and firms post unmet demand $\hat{X}_i = (X_i - Q_i)^+$ on the platform. (c) Given realized labor supply and demand posted on the platform, the platform operator carries out the allocation

¹⁴In practice, platforms sometimes quote this transfer payment in two parts: compensation to the temporary workers plus a commission fee. For ease of exposition, we do not model the commission fee explicitly, but one can easily find the corresponding commission cost by subtracting workers' wage from the transfer payment.

according to the announced policy. Then the understaffing cost $p((X_i - Q_i)^+ - A_i)^+$ occurs to each firm, and firms make transfer payments $\mathbf{w}(\mathbf{A}, \hat{\mathbf{X}})$ to the platform operator. Notationwise, we will use \hat{X}_i to simplify expressions as much as possible, but we will expand it as $(X_i - Q_i)^+$ if we want to emphasize its dependence on Q_i .

Remark 1. We make the following remarks on the sequence of events and modeling assumptions.

- (1) The platform operator does not need to observe each firm's permanent staffing level, and each firm does not need to observe other firms' staffing levels or the platform operator's decisions. In effect, the firms' staffing decisions and the platform operator's decision on the allocation policy and compensation rate constitute a simultaneous-move game. However, our proposed coordination mechanism will be valid even when either the platform or the firms move first.
- (2) In the above model, we have assumed that firms' vacancies $\hat{X}_i = (X_i - Q_i)^+$ are verifiable. In practice, this assumption is often satisfied, because the assigned workers can eventually verify the actual working time. If the actual working time is different from the requested time, the platform may intervene.¹⁵ Nevertheless, we will relax this assumption in Section 7 and show that the proposed mechanism can be adapted to ensure that firms will truthfully report their job vacancies.
- (3) In line with the supply chain coordination literature, in the analysis we consider the payment scheme $\mathbf{w}(\mathbf{A}, \hat{\mathbf{X}})$ to be determined by a central planner. As we show later, both the platform operator and individual firms will be willing to adopt the proposed payment scheme, as it can flexibly distribute the benefit of coordination among all players and make every player better off.

In what follows, we will first characterize the system-wide optimal solution $(\eta^*, \mathbf{A}^*(\cdot), \mathbf{Q}^*)$ which minimizes the total expected cost of the staffing system, and then propose a payment scheme $\mathbf{w}(\cdot)$ such that under the payment scheme the system-wide optimal solution $(\eta^*, \mathbf{A}^*(\cdot), \mathbf{Q}^*)$ constitutes a Nash equilibrium in the game the platform operator and n firms. Formally, $(\eta^*, \mathbf{A}^*(\cdot), \mathbf{Q}^*)$ is a Nash equilibrium if

- (i) given \mathbf{Q}^* , $(\eta^*, \mathbf{A}^*) \in \arg \max_{\eta, \mathbf{A}} \Pi(\eta, \mathbf{A}(\cdot) | \mathbf{Q}^*)$ as defined in problem (2); and
- (ii) for all $i \in I$, given $(\eta^*, \mathbf{A}^*(\cdot))$ and \mathbf{Q}_{-i}^* , $Q_i^* \in \arg \min_{Q_i \geq 0} h_i(Q_i | \mathbf{Q}_{-i}^*, \mathbf{A}^*(\cdot), \eta^*)$ as defined in problem (3).

¹⁵For example, in case the actual working time is much shorter than previously estimated, InStaff requires that employers pay 80% of the previously estimated amount. See <https://en.instaff.jobs/terms>.

That is, given firms' staffing levels fixed at \mathbf{Q}^* , the system-wide optimal compensation rate η^* and allocation rule $\mathbf{A}^*(\cdot)$ maximize the platform's payoff; given the platform choosing the compensation rate η^* and allocation rule $\mathbf{A}^*(\cdot)$ and other firms choosing staffing level \mathbf{Q}_{-i}^* , Q_i^* minimizes each Firm i 's staffing cost. Moreover, the proposed payment scheme makes all the parties better off with on-demand staffing. For readability, we will discuss the main ideas and sketch the proofs of some key results in the main text, but refer interested readers to Appendix A for the details.

4 System-wide Optimal Staffing Strategies: When Is an On-demand Workforce Beneficial?

In this section, we investigate when, from a system-wide point of view, an on-demand workforce can offer cost-saving benefits compared with a traditional staffing solution. This is equivalent to analyzing how the benefits of an on-demand workforce can be maximized when there are no conflicts of incentives from different parties in a centralized system. The results of this section will serve as a benchmark when we analyze the original problem, i.e., the decentralized system, in the next section.

From a system-wide perspective, the transfer payments are irrelevant in this analysis. For any feasible allocation policy, the system-wide optimization problem over $\mathbf{Q} = (Q_1, Q_2, \dots, Q_n)$, η , and $\mathbf{A}(\cdot)$ can be written as

$$\begin{aligned} \min_{\mathbf{Q}, \eta, \mathbf{A}(\cdot)} \quad & p \sum_{i \in I} \mathbb{E} [(X_i - Q_i)^+ - A_i((\mathbf{X} - \mathbf{Q})^+, S)] + \eta NG(\eta) + c \sum_{i \in I} Q_i \\ \text{s.t.} \quad & \mathbf{A}(\cdot) \in \mathcal{A} \\ & S = NG(\eta) \\ & Q_i \geq 0, \forall i \in I \end{aligned}$$

It is straightforward to see that the following no-waste condition holds in the optimal solution; otherwise, one can always reduce the understaffing cost by filling more job vacancies:

$$\sum_{i \in I} A_i(\hat{\mathbf{X}}, S) = \min \left\{ S, \sum_{i \in I} (X_i - Q_i)^+ \right\}, \text{ almost surely, } \forall \mathbf{X} \in \Omega. \quad (4)$$

Together with the relation $\min\{x, y\} = y - (y - x)^+$, we can simplify the objective function as

$$C(\mathbf{Q}, \eta) := p \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i)^+ - S \right)^+ \right] + \eta NG(\eta) + c \sum_{i \in I} Q_i. \quad (5)$$

Thus, the allocation policy $\mathbf{A}(\cdot)$ is irrelevant to the system-wide optimization. However, the allocation policy matters when the system is decentralized, and it is one of the key challenges in coordinating the system.

To characterize the structure of the solution in the general setting, we first consider two special cases in which only the permanent or on-demand workforce is available, and then derive the optimal staffing strategy when there is a mix of workforce. First, suppose that firms have no access to an on-demand staffing platform, and thus must rely solely on the permanent workforce. Minimizing (5) is reduced to $\min_{\mathbf{Q} \geq 0} \sum_{i \in I} \{p \mathbb{E} [(X_i - Q_i)^+] + cQ_i\}$. It is easy to see that the optimal permanent staffing level without on-demand workforce, denoted by Q_i^p , is determined by the first-order condition $p\bar{F}_i(Q_i^p) = c$. Second, suppose instead that the system can only use the on-demand workforce via the platform. By the log-concavity of G , it can be shown that the system-wide optimal compensation rate, denoted by η^o , is determined by

$$1 - F_{\sum X}(NG(\eta^o)) = \frac{1}{p} \left(\eta^o + \frac{G(\eta^o)}{g(\eta^o)} \right), \quad (6)$$

where $F_{\sum X}$ represents the cumulative distribution function of $\sum_{i \in I} X_i$.¹⁶

The following theorem shows that the system-wide optimal staffing strategy can be characterized with two thresholds ϕ^o and ϕ^p defined as $\phi^o = \eta^o + G(\eta^o)/g(\eta^o)$, and $\phi^p = p(1 - \mathbb{P}(X_i \leq Q_i^p, \forall i \in I)) = p\mathbb{P}(X_i > Q_i^p, \text{ for some } i \in I)$. In general, the optimal staffing strategy can use either the on-demand or permanent workforce exclusively, or a combination of both.

Theorem 1. (SYSTEM-WIDE OPTIMALITY) (i) If $\underline{\tau} \geq \phi^p$, it is optimal to exclusively use permanent workers: $\eta^* = \underline{\tau}$, and $Q_i^* = Q_i^p$ for all $i \in I$ where $Q_i^p = F_i^{-1}(p - c/p)$.

(ii) If $c \geq \phi^o$, it is optimal to exclusively use the on-demand workforce: $Q_i^* = 0$ for all $i \in I$, and $\eta^* = \eta^o$ as defined in (6).

(iii) Otherwise, mixed use of the permanent and on-demand workforces is optimal; optimal staffing levels and compensation rates satisfy the following optimality conditions:

$$p\mathbb{P} \left(NG(\eta^*) < \sum_{i \in I} (X_i - Q_i^*)^+ \right) = \eta^* + \frac{G(\eta^*)}{g(\eta^*)},$$

$$p\mathbb{P} \left(Q_i^* < X_i, NG(\eta^*) < \sum_{i \in I} (X_i - Q_i^*)^+ \right) = c, \forall i \in I.$$

Theorem 1 shows that the on-demand workforce complements permanent employees unless the workers' minimum availability threshold is too high, i.e., $\underline{\tau} \geq \phi^p$. Specifically, the system

¹⁶Because we have assumed that N is sufficiently large and $p > \underline{\tau}$, there exists a unique $\eta^o \in (\underline{\tau}, \bar{\tau})$ that satisfies equation (6).

should use the on-demand workforce exclusively if the permanent staffing cost c exceeds ϕ^o . Due to the self-scheduling behavior of on-demand workers, a higher compensation rate is required to attract more of them. Therefore, a mixed workforce can be optimal to balance the use of on-demand workforce with some permanent employees.

Whether the on-demand workforce benefits the system depends on the critical threshold ϕ_p . Analogous to the wisdom of inventory pooling (Eppen, 1979), an important benefit of the on-demand workforce lies in the risk-pooling effect as multiple random demands are filled by a common pool of on-demand workers. Therefore, the on-demand workforce would be more favorable, as the risk-pooling effect is stronger, as formalized below. Let \mathbf{X} and \mathbf{X}' be two n -dimensional random vectors. Then, \mathbf{X} is said to be smaller than \mathbf{X}' under the *supermodular order*, denoted by $\mathbf{X} \leq_{SM} \mathbf{X}'$, if $\mathbb{E}[\psi(\mathbf{X})] \leq \mathbb{E}[\psi(\mathbf{X}')] for all supermodular functions $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$. The supermodular order is a multivariate positive dependence order, which compares distributions exclusively based on their dependence (see, e.g., Corbett and Rajaram, 2006; Mak and Shen, 2014, and references therein). For example, let \mathbf{X} and \mathbf{X}' have the same (but arbitrary) marginal distributions with Normal copulae C and C' characterized by covariance matrices $\Sigma_{\mathbf{X}}$ and $\Sigma_{\mathbf{X}'}$, respectively. If $\Sigma_{\mathbf{X}} \leq \Sigma_{\mathbf{X}'}$ componentwise, that is, \mathbf{X}' has a higher pairwise dependence than \mathbf{X} , then $\mathbf{X} \leq_{SM} \mathbf{X}'$ (see Proposition 5 and its proof in Corbett and Rajaram, 2006).$

Corollary 1. (i) Suppose that X_i 's are independent across i . Then ϕ^p is increasing in n , and $\phi^p \rightarrow p$ as $n \rightarrow \infty$. That is, on-demand staffing always survives when there are sufficiently many employers with independent demands. (ii) For any two demand vectors \mathbf{X} and \mathbf{X}' such that \mathbf{X} is smaller than \mathbf{X}' in the supermodular order, written as $\mathbf{X} \leq_{SM} \mathbf{X}'$, the corresponding thresholds satisfy $\phi^p(\mathbf{X}) \geq \phi^p(\mathbf{X}')$.

The on-demand workforce will be always beneficial if the number of employers is sufficiently large and their demands are independent. This is because as n goes to infinity, the probability $\mathbb{P}(X_i > Q_i^p, \text{ for some } i \in I)$ will tend to one such that $\phi^p \approx p > \underline{\tau}$. Moreover, the on-demand workforce would be more attractive as firms' demands are less dependent across each other.

Corollary 2. $Q_i^* \leq Q_i^p$ for all $i \in I$.

Not surprisingly, with access to the on-demand workforce, firms should hire (weakly) fewer permanent employees than they would with the traditional staffing mode from a system-wide perspective. Corollary 2 highlights the importance of attracting individual firms to use the on-demand workforce with certain mechanisms. Self-interested firms may refuse to risk using the on-demand workforce unless a proper allocation of the on-demand workforce can be ensured and also aligned with the platform operator's incentive.

5 Coordination Mechanisms for the Decentralized System: Maximizing the Benefits of an On-Demand Workforce

We now consider a decentralized system in which the platform operator and n firms are self-interested. To exclude the uninteresting case in which the on-demand workforce should never be used (i.e., Part (i) of Theorem 1), we will henceforth assume $\underline{\tau} < \phi^p$. In such a system, the platform operator has to determine both the allocation policy and service contracts, and simultaneously, the firms decide their staffing strategies. The problem is challenging because the coordination here relies on the allocation policy and the payment scheme, $\mathbf{A}(\hat{\mathbf{X}}, S)$ and $\mathbf{w}(\mathbf{A}, \hat{\mathbf{X}})$, both of which can distort each firm's incentive. The allocation policy affects the labor supply that each firm is expected to get from the platform, and the payment scheme affects the cost of using the on-demand workforce. The two mechanisms need to be designed in the way that collectively, they induce all the firms to choose the optimal staffing strategies as in the centralized system to maximize the benefits of the on-demand workforce. To tackle this problem, we first discuss the allocation policy and then turn to the service contract design. We will show that our proposed mechanisms coordinate the decentralized system to achieve the first-best solution.

5.1 Allocation through Target Fill Rates

Although the allocation policy is irrelevant to the system-wide optimization, it will crucially influence how individual firms determine their optimal staffing levels in the decentralized system. In the existing literature on capacity pooling and allocation, a widely studied policy is called the relaxed linear allocation rule (see, e.g., Cachon and Lariviere, 1999; Netessine and Rudi, 2006), defined as

$$A_i^{lin}((\mathbf{X} - \mathbf{Q})^+, S) = \min \left\{ (X_i - Q_i)^+, (X_i - Q_i)^+ - \frac{1}{n} \left(\sum_{k=1}^n (X_k - Q_k)^+ - S \right) \right\} \forall i \in I. \quad (7)$$

Under the above allocation rule (and many others discussed in the literature), a firm's expected allocation depends on other players' staffing levels, leading to a noncooperative game among all the players. In Appendix B, we analyze our staffing system under the relaxed linear allocation rule along with a price-only contract—a typical setting considered in the literature. We show that due to the gaming effect, the equilibrium and system efficiency can substantially deviate from the system-wide optimality.

To resolve the above challenges, the ideal allocation policy should make each Firm i 's expected allocation $\mathbb{E}[A_i]$ independent of the other firms' job vacancies or permanent staffing levels. In what follows, we will first show that a novel fill rates-based allocation policy $\mathbf{A}^*(\cdot)$ can resolve the

above gaming effect, and then propose a form of contracts with target fill rates which induces the system-wide optimal solution in the decentralized system (and also incentivizes the self-interested platform operator to implement the proposed allocation policy).

By drawing the similarities between the inventory allocation in a pooling system and our on-demand workforce allocation, we leverage the Randomized Priority List policy developed in Zhong et al. (2017) based on online optimization to address the challenges in our context. For completeness, we summarize the allocation policy and describe it in our context as Algorithm 1 below. We name the policy as Target Fill Rate-Based (TFRB) policy to highlight the importance of target fill rates in coordinating the on-demand workforce system. Note that the policy has to be announced before the demand realization, and the firms can expect the received workforce levels based on their demand distribution and the announced policy. With slight abuse of notation, we use $\hat{\mathbf{X}}$ to denote the realized demand requests from the firms, and Algorithm 1 derives the amounts of workforce allocated to each firm. The policy takes in the target fill rate for each firm, denoted as $\beta_i \in [0, 1]$. We will show in Lemma 1 that if the allocation follows Algorithm 1, the target fill rates for all the firms can be satisfied in expectation and that the system-wide optimality can be achieved when the set of target fill rates are optimally chosen by solving a simple linear program (LP).

Lemma 1. The system achieves the minimum staffing cost $C(\mathbf{Q}^*, \eta^*)$ if the n firms set their permanent staffing levels as \mathbf{Q}^* and the platform operator offers a wage rate η^* and allocate the available on-demand workforce according to the TFRB policy $A^*(\cdot)$ with a set of target fill rates β^* , where \mathbf{Q}^* and η^* are as described in Theorem 1 and $\beta^* = (\beta_1^*, \beta_2^*, \dots, \beta_n^*)$ is an optimal solution to the linear program below:

$$\begin{aligned} \max_{\beta} \quad & \sum_{i \in I} \beta_i \mathbb{E}[(X_i - Q_i^*)^+] \\ \text{s.t.} \quad & \sum_{i \in U} \beta_i \mathbb{E}[(X_i - Q_i^*)^+] \leq \mathbb{E} \left[\min \left\{ S^*, \sum_{i \in U} (X_i - Q_i^*)^+ \right\} \right], \forall U \subseteq I \\ & \beta_i \geq 0, \forall i \in I \end{aligned} \quad (8)$$

Moreover, the expected allocation satisfies $\mathbb{E}[A_i^*((\mathbf{X} - \mathbf{Q}^*)^+, S^*)] = \beta_i^* \mathbb{E}[(X_i - Q_i^*)^+]$, for all $i \in I$.

The first set of constraints in problem (8) describes the feasibility conditions of the system. Given the workforce level S and the demand profiles, there are upper limits on the fill rates that can be satisfied through the allocation policies, which are captured by this set of inequalities. The last part of the above lemma shows that the expected amount of on-demand workers received

Algorithm 1 Target Fill Rate-Based (TFRB) Policy

* *Input:* Realization of requested demand $\hat{\mathbf{X}}$; distribution of $\hat{\mathbf{X}}$; workforce level S ; target service level $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_n)$.

1. Simulate T samples of realized requested demand using the distribution of $\hat{\mathbf{X}}$, where T is a large number. This gives rise to T realized requested demand vectors for the firms, denoted as $\hat{\mathbf{X}}(1), \hat{\mathbf{X}}(2), \dots, \hat{\mathbf{X}}(T)$. Treat these as the requested demands for T periods, i.e., $\hat{\mathbf{X}}(t)$ is the requested demand in simulated period t for $t = 1, 2, \dots, T$. Assume the allocation problem is repeated T times with S amount of workforce to allocate in each period.

a. Initiate $t = 1$. Allocate the workforce arbitrarily for requested demand $\hat{\mathbf{X}}(1)$. The initial allocation will not affect the convergence results of the policy. Set $t = 2$.

b. In the simulated period t , for $t > 1$, allocate the workforce, S , based on a priority rule that firms with higher priorities will be satisfied fully before those with lower priorities. The allocation ends either when the workforce is depleted or all the requested demands are satisfied. The priority of Firm i is based on the gap between the expected allocation in the past $(t - 1)$ simulated periods and the actual allocation received, i.e., $(t - 1)\beta_i \mathbb{E}[\hat{X}_i] - \sum_{s=1}^{t-1} A_i(s)$, where with a little abuse of notation, $A_i(s)$ represents the actual workforce allocated to Firm i in simulated period s . Firms with larger gaps in the past have higher priorities in the current period. Let $[i]$ denote the index of the firm at the i^{th} position in the priority list. There are two possible scenarios:

- *Scenario 1:* $\sum_{i \in I} \hat{X}_{[i]}(t) > S$, i.e., the workforce is not sufficient to satisfy the requested demands of all firms. Define $i' = \min\{j : \sum_{i=1}^j \hat{X}_{[i]}(t) > S\}$. Then the allocation stops at Firm $[i']$ with $A_{[i]}(t) = \hat{X}_{[i]}(t)$, for all $i < i'$, and Firm $[i']$ receives the remaining amount of workforce, i.e., $A_{[i']}(t) = S - \sum_{j=1}^{i'-1} \hat{X}_{[j]}(t)$. The remaining firms receive zero allocation.
- *Scenario 2:* $\sum_{i \in I} \hat{X}_{[i]}(t) \leq S$, i.e., the workforce is sufficient to satisfy the requested demands of all firms. Then $A_{[i]}(t) = \hat{X}_{[i]}(t)$, for all $i \in I$.

Set $t = t + 1$.

c. Repeat Step b until $t > T$. Denote the priority list used in simulated period t as $L(t)$, for $t > 1$. By this step, we have generated $(T - 1)$ priority lists, i.e., $L(2), L(3), \dots, L(T)$. Each priority list corresponds to a ranking of n firms.

2. Randomly draw one priority list from the set of $(T - 1)$ priority lists generated in Step 1, $\{L(2), L(3), \dots, L(T)\}$.

3. Allocate \hat{X}_i amount of workforce to Firm i following the priority list obtained in Step 2 until reaching the end of the list, or the available workforce is depleted.

* *Output:* An allocation of on-demand workforce $\mathbf{A}(\hat{\mathbf{X}}, S)$ obtained from Step 3.

by Firm i is independent of the other firms' requests if the optimal target fill rate β_i^* is specified in a contract beforehand. This can therefore resolve the gaming effect discussed earlier.

We note that there are generally multiple solutions of β^* to the LP problem (8). Under a corner solution of the LP, some firms may be offered with much lower fill rates than others. To refine the solution such that the offered fill rates are close to each other as much as possible, we can follow the procedure below.

Algorithm 2 Refinement of Target Fill Rates

* *Input:* System-wide optimal S^* and \mathbf{Q}^* and distribution of \mathbf{X} to evaluate $\mathbb{E}[(X_i - Q_i^*)^+]$ for all $i \in I$ and $\mathbb{E}[\min\{S^*, \sum_{i \in U}(X_i - Q_i^*)^+\}]$ for all $U \subseteq I$.

1. Solve the following LP to maximize the lower bound of all feasible fill rates:

$$\begin{aligned} \max_{\beta, b} \quad & b \\ \text{s.t.} \quad & \sum_{i \in U} \beta_i \mathbb{E}[(X_i - Q_i^*)^+] \leq \mathbb{E} \left[\min \left\{ S^*, \sum_{i \in U} (X_i - Q_i^*)^+ \right\} \right], \forall U \subseteq I \\ & \beta_i \geq b, \forall i \in I. \end{aligned} \quad (9)$$

Obtain b^* as the optimal objective value of problem (9).

2. Solve problem (8) with an additional set of constraints:

$$\beta_i \geq b^*, \forall i \in I. \quad (10)$$

* *Output:* Target fill rates $\beta^* =$ the optimal solution from Step 2.

In Algorithm 2, we first find the maximum lower bound of all feasible fill rates by solving problem (9). Then, using the maximum lower bound b^* , we refine the feasible set of the original problem (8) by requiring all the β_i 's not less than b^* . Thus, solving the refined LP will help avoid assigning an extremely low fill rate to any firm.¹⁷ If the original LP has an optimal solution such that all the β_i^* 's are equal (which may not be true in general), then we have $\beta_i^* = b^*$ for all $i \in I$ and Algorithm 2 will identify that solution. Our extensive numerical study suggests that using Algorithm 2, one can significantly reduce the difference in the optimal target fill rates associated with each firm. Moreover, in 83% of the instances we tested (including ones with highly asymmetric firms), the refined fill rates by Algorithm 2 are identical¹⁸ to all firms. Detailed numerical results are reported in Appendix C.

¹⁷We are grateful to the anonymous Associate Editor for suggesting this approach.

¹⁸The refined fill rates are said to be identical if we observed the difference between the maximum and minimum fill rates offered is less than 10^{-6} .

5.2 Coordination Contract

We first propose a form of contracts based on the optimal target fill rates β^* and then verify that it coordinates the system when contract parameters are appropriately chosen. We say a contract $w(\cdot)$ coordinates the decentralized system if (i) under the proposed contract, the system-wide optimal permanent staffing levels \mathbf{Q}^* and compensation rate η^* along with the TFRB policy $\mathbf{A}^*(\cdot)$ constitute a Nash equilibrium defined in Section 3, and (ii) the proposed contracts provide the platform operator with a nonnegative expected profit and not render individual firms worse off by joining the platform.

We consider the following form of service contracts between the platform operator and each Firm i , which consists of a membership fee r_i , a fixed rate w_i^F and a contingent payment $w_i^C(A_i, \hat{X}_i)$ based on the target fill rate β_i^* . Note that the values of β_i^* 's can be generated according to Lemma 1 beforehand.

- **Membership Fee.** To enjoy the platform's on-demand staffing services, each firm pays a membership fee for registration, denoted by r_i , to the platform operator.

For the other contract terms, we introduce a set of parameters m_i 's; as will be shown later, the value of each m_i can be properly chosen to coordinate the system.

- **Fixed Rate.** For each unit of an on-demand workforce allocated, Firm i pays a fixed rate to the platform

$$w_i^F = p \left(1 - \frac{m_i}{\beta_i^*} \right). \quad (11)$$

- **Contingent Payment.** Depending on the actual allocation A_i and job vacancy $\hat{X}_i = (X_i - Q_i)^+$, Firm i makes a contingent payment defined as

$$w_i^C(A_i, \hat{X}_i) = \frac{pm_i(A_i - \beta_i^* \hat{X}_i)}{\beta_i^*}. \quad (12)$$

The above contingent payment $w_i^C(A_i, \hat{X}_i)$ has an intuitive interpretation. If the actual allocation oversatisfies the fill-rate based target, i.e., $A_i > \beta_i^* \hat{X}_i$, the term $w_i^C(A_i, \hat{X}_i)$ serves as a price premium. If the delivered allocation is below the target, i.e., $A_i < \beta_i^* \hat{X}_i$, then $w_i^C(A_i, \hat{X}_i)$ is negative and can be viewed as a subsidy to Firm i . From the practical perspective, the proposed contingent payment is in line with the concept of performance-based contracting used in the on-demand staffing industry, among many other service industries, in that it separates the customer's expectations of service and the service provider's actual performance (Kim et al., 2007). Based on our interview with a leading staffing platform in the Netherlands that provides

on-demand workers to, e.g., warehouses owned by retailers and logistics providers, performance-based contracts are commonly used. For example, its service contracts with clients are contingent on the actual number of picks performed at a warehouse relative to a target number with each client; a bonus or minus payment can be applied depending on whether the delivered number is higher or lower than the target.¹⁹ Another example is call center staffing services, for which the contract terms often depend on the number of calls answered compared to a target number specified in the service agreement (Robbins and Harrison, 2011).

5.2.1 Platform Operator

We first consider the platform operator's problem under the proposed contract. The platform operator must pay $\eta S = \eta NG(\eta)$ to compensate the on-demand workers available and deploy them in an on-demand fashion. She will thus incur an overstaffing cost whenever the number of available on-demand workers exceeds aggregate demand $\sum_{i \in I} \hat{X}_i$. Therefore, the platform operator may target a lower on-demand staffing level than the system-wide optimal level S^* unless the payment scheme penalizes her appropriately in the event of a staffing shortfall. As will be shown below, the contingent term $w_i^C(A_i, \hat{X}_i)$ exactly plays this role in balancing the platform operator's overage and underage risks.

Given individual firms' permanent staffing levels fixed as \mathbf{Q} , under the proposed payment scheme the platform operator's problem (2) can be reduced to

$$\begin{aligned} \max_{\eta, \mathbf{A}(\cdot)} \Pi(\eta, \mathbf{A}(\cdot) | \mathbf{Q}) &= -\eta NG(\eta) + \sum_{i \in I} \mathbb{E} \left[w_i^F A_i + w_i^C(A_i, \hat{X}_i) + r_i \right] \\ &= -\eta NG(\eta) + p \sum_{i \in I} \mathbb{E}[A_i] - p \sum_{i \in I} m_i \mathbb{E}[\hat{X}_i] + \sum_{i \in I} r_i \\ &= -\eta NG(\eta) + p \mathbb{E} \left[\min \left\{ S, \sum_{i \in I} \hat{X}_i \right\} \right] - p \sum_{i \in I} m_i \mathbb{E}[\hat{X}_i] + \sum_{i \in I} r_i, \end{aligned}$$

where $S = NG(\eta)$. The first equality follows by the definitions of w_i^F and $w_i^C(\cdot)$. The second equality follows by invoking the no-waste condition (4) for the allocation policy. We can further

¹⁹In our model, $\beta_i^* \hat{X}_i$ and A_i represent the amounts of the workforce (e.g., total working hours by the on-demand workforce) targeted and delivered; they can be readily converted to some productivity measures, e.g., the number of picks, by multiplying the total hours by the average productivity per worker per hour.

rewrite the platform operator's profit as

$$\begin{aligned}
\Pi(\eta, \mathbf{A}(\cdot)|\mathbf{Q}) &= -\eta NG(\eta) - c \sum_{i \in I} Q_i - p \mathbb{E} \left[\left(\sum_{i \in I} \hat{X}_i - S \right)^+ \right] + c \sum_{i \in I} Q_i + p \mathbb{E} \left[\left(\sum_{i \in I} \hat{X}_i - S \right)^+ \right] \\
&\quad + p \mathbb{E} \left[\min \left\{ S, \sum_{i \in I} \hat{X}_i \right\} \right] - p \sum_{i \in I} m_i \mathbb{E} [\hat{X}_i] + \sum_{i \in I} r_i \\
&= -C(\mathbf{Q}, \eta) + p \sum_{i \in I} (1 - m_i) \mathbb{E} [(X_i - Q_i)^+] + c \sum_{i \in I} Q_i + \sum_{i \in I} r_i,
\end{aligned}$$

where the last equality follows by the definition of $C(\mathbf{Q}, \eta)$, i.e., the system's cost function as defined in (5), the relation $(x - y)^+ + \min\{y, x\} = x$, and $\hat{X}_i = (X_i - Q_i)^+$. Thus, maximizing Π over η is equivalent to minimizing $C(\mathbf{Q}, \eta)$ when permanent staffing levels are being fixed as \mathbf{Q} . Therefore, given $\mathbf{Q} = \mathbf{Q}^*$, the platform operator will set η to be the system-wide optimal value η^* in the same manner as we did in Section 4.

Furthermore, it is easy to see that the platform operator is indifferent among all allocation policies that satisfy the no-waste condition, and so it will be optimal for the platform to assign the on-demand workforce according to the TFRB policy.²⁰

Lemma 2 below summarizes the above discussion and concludes that the proposed contract aligns the platform operator's incentive with the system.

Lemma 2. Suppose that the permanent staffing levels are fixed as \mathbf{Q}^* . The payment schemes w_i^F and w_i^C , defined in (11) and (12), incentivize the platform operator to choose the system-wide optimal compensation rate η^* and to voluntarily allocate the on-demand workforce via the TFRB policy $\mathbf{A}^*(\cdot)$ such that the optimal target fill rates β^* as characterized in Lemma 1 are warranted in expectation.

5.2.2 Individual Employers

Next, we turn to the incentive issue for individual firms. The decentralized system can be coordinated only if the system-wide optimal permanent staffing levels \mathbf{Q}^* are chosen by self-interested firms. We should ensure that under the proposed contracts, Q_i^* minimizes $h_i(Q_i)$, as defined in (3) for each Firm i , given that the platform operator chooses the system-wide optimal compensation rate η^* and follows the TFRB policy $\mathbf{A}^*(\cdot)$.

Generally, firms may engage in a noncooperative game when choosing \mathbf{Q} . The proposed mechanism eliminates this issue and dramatically simplifies each firm's problem. Under the

²⁰Here, we invoke the conventional assumption that when facing a set of indifferent allocation policies, the platform operator will choose the one that favors the system. In our setting, this can be justified because the platform operator may encourage firms' participation and accrue goodwill by following the desired allocation policy.

contract mechanism, Firm i 's expected payment to the platform is $\mathbb{E}[w_i^F A_i + w_i^C(A_i, \hat{X}_i) + r_i] = p\mathbb{E}[A_i - m_i(X_i - Q_i)^+] + r_i$, and its expected understaffing cost is $p\mathbb{E}[(X_i - Q_i)^+ - A_i]$. Thus, Firm i 's problem can be reduced to

$$\min_{Q_i \geq 0} h_i(Q_i) = p(1 - m_i)\mathbb{E}[(X_i - Q_i)^+] + cQ_i + r_i. \quad (13)$$

The formulation highlights the critical impact of contract parameter m_i on Firm i 's permanent staffing level. The following lemma demonstrates that with a properly chosen m_i , Firm i 's optimal permanent staffing level aligns with the system-wide optimal value of Q_i^* .

Lemma 3. Under the contract (r_i, w_i^F, w_i^C) , each Firm i will choose the system-wide permanent staffing level Q_i^* if parameters m_i 's are chosen such that

$$m_i = \frac{p\bar{F}_i(Q_i^*) - c}{p\bar{F}_i(Q_i^*)}, \forall i \in I. \quad (14)$$

Note that m_i is decreasing in Q_i^* —or, equivalently, increasing in the expected usage of on-demand workforce—which indicates that to achieve coordination, the platform should share more of the cost, as the on-demand workforce contributes more to the system.

5.2.3 Individual Participation

It remains to address the participation issue. Under our proposed allocation and contract mechanism, the platform operator receives an expected profit Π^* :

$$\Pi^* = p \sum_{i \in I} (1 - m_i)\mathbb{E}[(X_i - Q_i^*)^+] - C(Q^*, \eta^*) + \sum_{i \in I} r_i, \quad (15)$$

where the values of m_i 's are set according to (14). As the values of m_i 's must be large enough to incentivize employers, the first two terms can be negative. Nevertheless, the platform operator can leverage membership fees to collect a positive surplus. Clearly, membership fees r_i 's should be chosen such that

$$\sum_{i \in I} r_i \geq C(Q^*, \eta^*) - p \sum_{i \in I} (1 - m_i)\mathbb{E}[(X_i - Q_i^*)^+]. \quad (16)$$

On the other hand, let $h_i^p = p\mathbb{E}[(X_i - Q_i^p)^+] + cQ_i^p$, representing the staffing cost incurred by Firm i when he declines to join the platform. From the discussion in Section 5.2.2, it is straightforward to see that to ensure each firm's participation, the membership fee r_i should be

less than Firm i 's cost reduction by on-demand staffing, i.e.,

$$r_i \leq h_i^p - \{p(1 - m_i) \mathbb{E} [(X_i - Q_i^*)^+] + cQ_i^*\}. \quad (17)$$

Lastly, since the system attains optimality, the overall staffing cost is reduced by the on-demand workforce. Thus, there exists a set of membership fees r_i 's that satisfy (16) and (17).

Together with Lemmas 2 and 3, we have shown that under the proposed contract, the platform operator and the n individual firms will voluntarily choose the system-wide optimal solution in equilibrium and that all parties benefit from the on-demand workforce. Note that the proposed mechanism cannot lead to any other Nash equilibrium that is strictly worse than the system optimality. This is because by Lemma 3, Q_i^* is a dominant strategy of Firm i , and so the platform operator will always consider her problem with all $Q_i = Q_i^*$, which results in system optimality by Lemma 2. Therefore, we can conclude Section 5's main result in Theorem 2 below.

Theorem 2. (SYSTEM COORDINATION) The contracts (r_i, w_i^F, w_i^C) coordinate the system when parameters m_i 's are chosen according to (14) and membership fees r_i 's satisfy (16) and (17).

To recap the intuition behind our coordination contracts, note that summing up the three contract terms results in a total payment of $r_i + p(A_i - m_i \hat{X}_i)$, which is linear in the actual allocation and job vacancy. Broadly speaking, our problem can be viewed as a double-sided moral hazard problem (e.g., Bhattacharyya and Lafontaine, 1995) in the sense that the platform's decisions η and $\mathbf{A}(\cdot)$ and each firm's permanent staffing level Q_i jointly determine the system cost but none of these decisions is contractible.²¹ The term $p(A_i - m_i \hat{X}_i)$ essentially determines a risk-sharing rule between each firm and the platform operator. As job vacancy $\hat{X}_i = (X_i - Q_i)^+$ depends on each firm's decision on Q_i , a set of properly chosen m_i 's will incentivize firms to choose the system-wide optimal staffing levels \mathbf{Q}^* , which in turn aligns the platform operator's incentive with the system-wide optimality. Additionally, the membership fees in the contract not only guarantee individual participation but also enable the gain from coordination to be flexibly redistributed among all parties in the system. As r_i decreases, Firm i will enjoy more benefit from coordination, whereas the platform operator will surrender more surplus. Thus, the platform operator and individual firms can bargain over membership fees to divide the total surplus without affecting system efficiency.

Finally, we close this section by noting that the possible contract terms are affected by the benefit of inventory pooling. Corollary 3 below shows that the system can be coordinated without any membership fees if maintaining permanent workers is costly, and such membership-

²¹Our problem departs from the classic double-sided moral hazard framework (e.g., Bhattacharyya and Lafontaine, 1995) as it involves $(n + 1)$ players and the allocation of a common resource.

free contracts are more likely to achieve coordination as demands are less positively dependent (i.e., the risk pooling effect is stronger).

Corollary 3. Assuming that exclusively using the on-demand workforce is system-wide optimal, i.e., $Q^* = 0$, the system can be coordinated without membership fees, i.e., $r_i = 0$ for all $i \in I$ as long as $c \sum_{i \in I} \mathbb{E}[X_i] \geq C(0, \eta^o)$. Furthermore, $C(0, \eta^o)$ decreases as \mathbf{X} becomes smaller in the supermodular order.

6 Firm-Independent Contracts

6.1 Exchangeable Demands

We have so far allowed the contract terms to be firm-dependent. In some situations, it may be easier to implement a contract if the contract terms are identical to all firms. In fact, the firm-dependent feature of the proposed contract is due to that firms' demands may follow different marginal distributions, and their dependency may not be symmetric. We say random demands (X_1, X_2, \dots, X_n) are *exchangeable* if the joint distribution of (X_1, X_2, \dots, X_n) is the same as that of $(X_{(1)}, X_{(2)}, \dots, X_{(n)})$ for any permutation $\{(1), (2), \dots, (n)\}$ of $\{1, 2, \dots, n\}$. The exchangeability of (X_1, X_2, \dots, X_n) implies that the marginal distribution X_i is identical across firms and that the dependence structure of each firm's demand on the others' is the same. It is a weaker condition than requiring the X_i 's to be independent and identically distributed. In what follows, we establish that there will be a firm-independent contract that can coordinate the system as long as firms' demands are exchangeable.

By Theorem 1, under the exchangeability condition, the system-wide optimal staffing levels Q_i^* for each firm are the same. Let Q^* denote the system-wide staffing level of individual firms. The following proposition characterizes the firm-independent contract in closed form, provided the exchangeability condition. Notably, in such a contract, there is an identical target fill rate β^* applied to all firms.

Proposition 1. Assuming that random demands (X_1, X_2, \dots, X_n) are exchangeable with an identical marginal distribution $F(x)$, there exists a firm-independent contract (r, w^F, w^C) that coordinates the system. In the contract, we have $w^F = p(1 - m/\beta^*)$, $w^C(A_i, \hat{X}_i) = pm(A_i - \beta^* \hat{X}_i)/\beta^*$, and any membership fee r satisfying (16)–(17), where the firm-independent target fill rate is given by

$$\beta^* = \frac{\mathbb{E}[\min\{S^*, \sum_{i \in I} (X_i - Q^*)^+\}]}{n \int_{Q^*}^{\infty} (x - Q^*) dF(x)}. \quad (18)$$

and the contract parameter $m = 1 - c/(p\bar{F}(Q^*))$.

Note that the firm-independent target fill rate in (18) implies that the expected allocation of each Firm i is given by $\mathbb{E}[A_i^*] = \mathbb{E}[\min\{S^*, \sum_{i \in I} (X_i - Q_i^*)^+\}]/n$. That is, firms will equally likely share the on-demand workforce. This is because when all target fill rates are equal and demands are exchangeable, the priority list used in the TFRB policy will be equally likely drawn from all the firm permutations.

6.2 Nonexchangeable Demands

When demands are not exchangeable but the contract terms need to be firm-independent for practical reasons, it is worth exploring whether a firm-independent contract can achieve a reasonably good performance. We consider a firm-independent contract constructed in the following manner. Notice that the proposed contract terms w_i^F and $w_i^C(\cdot)$ defined in (11)–(12) depend on each specific firm's demand distribution only through the optimal target fill rate β_i^* and parameter m_i , where m_i as defined in (14) depends on $\bar{F}_i(Q_i^*)$, i.e., Firm i 's understaffing probability without on-demand staffing. To make contract terms firm-independent, (i) we use the best lower bound of feasible fill rates, i.e., b^* obtained by solving problem (9), instead of the exact optimal target fill rates β_i^* 's, to calculate payment terms and, (ii) set

$$m_i = \bar{m} = 1 - \frac{c}{p \left(\sum_{i \in I} \bar{F}_i(Q_i^*)/n \right)}, \forall i \in I. \quad (19)$$

That is, we approximate parameter m_i with the average probability $\sum_{i \in I} \bar{F}_i(Q_i^*)/n$. As such, the firm-independent contract terms are given by $w^F = p(1 - \bar{m}/b^*)$ and $w^C(A_i, \hat{X}_i) = p\bar{m}(A_i - b^*\hat{X}_i)/b^*$ along with any membership fee r satisfying (16)–(17). Note that b^* is a nominal fill rate used only for the calculation of the firm-independent contract terms. It can be different from the fill rates delivered by the platform in equilibrium, as explained below.

The proposed firm-independent contract has the same contractual form as the coordination contract analyzed in Section 5.2 except that parameter \bar{m} differs from m_i . Following the same approach as in Section 5.2, we can verify that two appealing properties of the original coordination contract are preserved: (1) The platform operator is incentivized to minimize the system cost $C(\mathbf{Q}, \eta)$ over η for any given \mathbf{Q} while being indifferent among all the no-waste allocation policies; (2) each Firm i 's problem reduces to a simple cost minimization over Q_i independent of the actions of other firms and the platform. The resulting equilibrium is summarized below.

Lemma 4. Under firm-independent contract terms, $w^F = p(1 - \bar{m}/b^*)$, $w^C(A_i, \hat{X}_i) = p\bar{m}(A_i - b^*\hat{X}_i)/b^*$, and any membership fee r satisfying (16)–(17), in equilibrium, each firm chooses a permanent staffing level Q_i^I that minimizes its staffing cost $h_i(Q_i) = p(1 - \bar{m})\mathbb{E}[(X_i - Q_i)^+] +$

cQ_i+r , whereas the platform operator sets a compensation rate η^I that minimizes the system cost given $Q_i = Q_i^I$ for all $i \in I$, and voluntarily adopts the TFRB policy to allocate the on-demand workforce. The equilibrium leads to a system cost $C(\mathbf{Q}^I, \eta^I)$.

The target fill rates of the TFRB policy in equilibrium, denoted by $\beta^I = (\beta_1^I, \dots, \beta_n^I)$, can be determined by Algorithm 2 but with the equilibrium staffing levels \mathbf{Q}^I and $S^I = NG(\eta^I)$ as input. Since every firm now has identical contract terms, the platform can randomize firm indexes so that each firm will be assigned with any target fill rate β_i^I from Algorithm 2 with equal probability. Then, $\bar{\beta}^I = \sum_{i=1}^n \beta_i^I/n$ will be the uniform fill rate expected to be delivered to all firms. Note that $\bar{\beta}^I$ or β_i^I 's are not necessarily equal to b^* in Algorithm 2, which is used only to calculate contract terms. A more detailed discussion on the equilibrium outcome under firm-independent contracts can be found in the proof of Lemma 4.

Because \bar{m} is not generally equal to m_i , the Q_i^I 's chosen by each firm can deviate from the system-wide optimal solution Q_i^* , so can the platform's compensation rate η^I . Note that the approximation \bar{m} is exact when $Q_i^* = 0$ for all $i \in I$. It follows immediately that the aforementioned firm-independent contract coordinates the system if the system-wide solution entails exclusive use of the on-demand workforce, as summarized in Proposition 2.

Proposition 2. Suppose that the on-demand workforce should be used exclusively in the system-wide optimal solution, i.e., $Q_i^* = 0$ for all $i \in I$. Then, a firm-independent contract (r, w^F, w^C) coordinates the system with $w^F = p(1 - \bar{m}/b^*)$, $w^C(A_i, \hat{X}_i) = p\bar{m}(A_i - b^*\hat{X}_i)/b^*$ and any membership fee r satisfying (16)–(17), where b^* is the optimal objective value to problem (9) and $\bar{m} = (p - c)/p$.

For cases where $Q_i^* > 0$ for some $i \in I$, however, the firm-independent contract could lead to some inefficiency. We define the loss of efficiency (due to restricting to the firm-independent contract) as $(C(\mathbf{Q}^I, \eta^I) - C(\mathbf{Q}^*, \eta^*)) / C(\mathbf{Q}^*, \eta^*)$. Below we numerically investigate the efficiency loss under different parameter settings.

We consider four firms ($n = 4$) with parameters $p = 10$, $N = 1000$, and $G \sim Uniform(3, 9)$. To simulate the asymmetric demand distributions, we set random demands as follows. Assume the mean of total demand to be fixed at 400. Let the mean demand of Firm 1 be 400γ , and the mean demands of other firms be $400(1 - \gamma)/3$, where $\gamma \in (0, 1)$. That is, γ represents the (average) demand share of Firm 1, while the mean demands of the other three firms are equal. For example, when $\gamma = 0.25$, all the mean demands are equal; when $\gamma = 0.9$ ($= 0.1$), Firm 1 represents an exceptionally big (small) customer of the platform. The marginal distributions of the X_i 's are chosen from Gamma, log-normal, and uniform distributions with coefficient of

variation (CV) in $\{0.5, 1, 1.5\}$. For uniformly distributed demand, we restrict to a CV of 0.5 to guarantee nonnegative support.

The joint distribution of (X_1, X_2, X_3, X_4) is constructed based on a Gaussian copula with correlation matrix Σ , where $\Sigma_{ii} = 1$ and $\Sigma_{ij} = \rho$ for all $i, j \in \{1, 2, 3, 4\}$ and $i \neq j$. Here, ρ captures the demand correlation across firms. For any Firm 1's demand share $\gamma \in \{0.05, 0.1, 0.4, 0.6, 0.9, 0.95\}$, any permanent staffing cost $c \in \{4, 5, 6\}$ and any correlation $\rho \in \{0, 0.5\}$, we evaluate the loss of efficiency in the following four sets of experiments: (1) all marginal distributions are Gamma with an equal CV varied in $\{0.5, 1, 1.5\}$; (2) all marginal distributions are log-normal with an equal CV varied in $\{0.5, 1, 1.5\}$; (3) all marginal distributions are uniform with an equal CV = 0.5; (4) the marginal distribution and CV of each X_i are randomly drawn from $\{\text{Gamma, log-normal, uniform}\}$ and $\{0.5, 1, 1.5\}$ respectively with equal probabilities. Note that in the fourth set of experiments, we test 10 randomly generated mixtures of marginal distributions for any parameter combinations. In total, there are 714 instances tested.

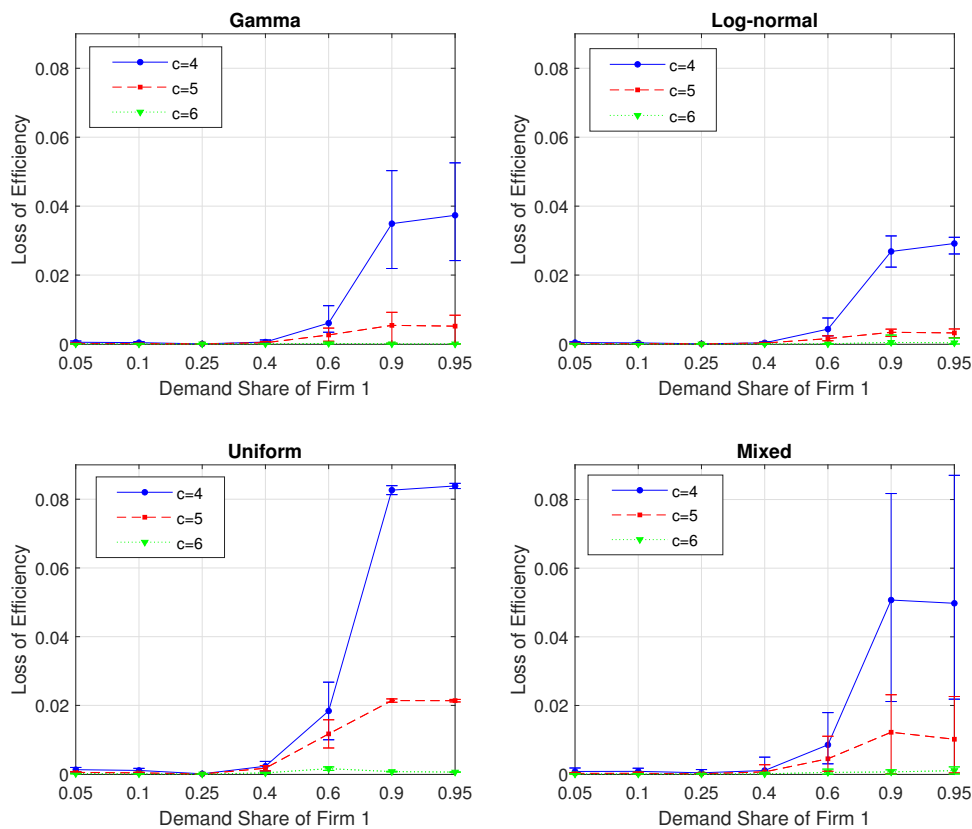


Figure 1: [Color online] Performance of firm-independent contracts

Our numerical results are summarized in Figure 1, where we plot the average loss of efficiency along with error bars showing its maximum and minimum. The firm-independent contract performs very well with a loss of efficiency less than 1% if $\gamma \leq 0.4$, i.e., Firm 1's demand share is

not too different from others (which echoes Proposition 1 we established earlier for exchangeable demands). Even when Firm 1's demand accounts for 60% of the total demand, in most instances, the loss of efficiency is below 2%. However, the loss of efficiency becomes larger if Firm 1's demand share is extremely large; the worst case we observed has a loss of efficiency of 8.70%, which occurs when Firm 1's demand accounts for 95% of the total demand and the marginal distributions are mixed (see the bottom-right panel of Figure 1). Moreover, we observe that the loss of efficiency is close to zero for $c = 6$ because with large permanent staffing cost c , the system-wide optimal solution uses the on-demand workforce almost exclusively, thus making the firm-independent contract near-optimal as proved in Proposition 2. We also note that for the instances with $c = 4$, most workers used in the system-wide optimal solution are permanent; thus, the selected range of c is wide enough to cover most possible scenarios, including ones where the on-demand workforce is not a primary source of labor.

In summary, our numerical study confirms that a firm-independent contract can induce near-optimal system performance when (1) firms' demands are fairly similar or (2) the on-demand workforce is significantly cost-efficient compared to the permanent workforce. Furthermore, customizing contracts for exceptionally large customers is valuable. Since w_i^F decreases with understaffing probability without an on-demand workforce, $\bar{F}_i(Q_i^*)$, firms needing more on-demand workers should get lower fixed rates in the coordinating contract. In practice, it is not uncommon to offer discounts to big customers.

7 Unverifiable Job Vacancies

In our basic model, we have assumed that job vacancies can be verified (e.g., by on-demand workers performing the job) so that firms would always report their actual demand to the platform. This is true for some but not all settings in reality. In this section, we show how the proposed mechanism in Section 5 can be adapted to elicit truthful demand information.

Suppose now that $\hat{X}_i = (X_i - Q_i)^+$ is Firm i 's private information. Firms may potentially lie about their actual demands to the platform. Note that it is the contingent term $w_i^C(A_i, \hat{X}_i)$ which induces firms to lie, since firms can receive extra compensation by exaggerating \hat{X}_i . Inspired by this observation, we propose introducing a third party to manage contingent payment $w_i^C(\cdot)$ such that each firm only pays a membership fee and the fixed rate for each unit of the on-demand workforce. The third party can be an independent industry body set up by the government or management services providers. The idea of having a third party to manage cash flows is reminiscent of Shang et al. (2009); see more examples of potential third parties therein.

We consider the subgame in which permanent staffing levels Q_i 's and the platform's com-

pensation rate η have been settled. The sequence of events is as follows: (1) X_i 's are realized, but only observed privately by each firm; (2) Firm i reports its job vacancies as \check{X}_i ; (3) the on-demand workforce is then allocated to each firm based on the TFRB rule, according to reported \check{X}_i ; and (4) payments are transferred under the following fixed-rate scheme.

- **Fixed-Rate Payment.** For each unit of allocated on-demand workforce, Firm i pays the platform operator a fixed rate $w_i^F = p(1 - m_i/\beta_i^*)$ along with a membership fee r_i ;
- **Third-Party Transfer.** Based on each actual allocation A_i and reported demand \check{X}_i , $t_i = r_0 - w_i^C(A_i, \check{X}_i)$ is transferred from the platform operator to a third party, where $r_0 \geq 0$.

Lemma 5. If realized job vacancies \hat{X}_i 's are not verifiable, the fixed-rate payment scheme induces firms to truthfully report their job vacancies, i.e., $\check{X}_i = \hat{X}_i$, for all $i \in I$ provided $w_i^F \geq 0$.

The lemma establishes that the fixed-rate payment scheme induces truth-telling. Intuitively, with a fixed-rate payment, when the on-demand workforce is sufficient, a firm incurs overstaffing costs if overreporting \hat{X}_i . When the on-demand workforce is insufficient, overreporting does not help increase the firm's allocation due to the TFRB policy. Hence, firms cannot benefit from exaggerating actual job vacancies. To ensure truthful-telling, we need a mild condition $w_i^F \geq 0$, or equivalently, $m_i \leq \beta_i^*$.²² This condition is satisfied in around 90% of the instances of our numerical study reported in Section 6. Moreover, a sufficient condition for $w_i^F \geq 0 \forall i$ is derived in Appendix D. For cases of $w_i^F < 0$, we can easily modify the above mechanism to ensure truth-telling at little efficiency loss, as discussed at the end of this section.

Notably, removing contingent term $w_i^C(A_i, \hat{X}_i)$ from each firm's payment will not affect their incentives in choosing the permanent staffing level, as the TFRB allocation rule guarantees $\mathbb{E}[w_i^C(A_i, \hat{X}_i)] = 0$. On the other hand, the third-party transfer ensures that the contingent term $w_i^C(\cdot)$ applies to the platform as before such that the platform operator's incentive is aligned with the system-wide optimality. As $\mathbb{E}[w_i^C(A_i, \hat{X}_i)] = 0$, the fixed amount r_0 in the third-party transfer is used to engage the third party by providing it with a nonnegative expected payoff. We therefore have the following proposition.

Proposition 3. If job vacancies are unverifiable, provided $w_i^F \geq 0$ (or equivalently, $m_i \leq \beta_i^*$) for all $i \in I$, the contracts (r_i, w_i^F) , i.e., a membership fee plus a fixed-rate payment, coordinate

²²In general, w_i^F could be negative if a firm needs to be provided with a very strong incentive to use the on-demand workforce. In Appendix D we have proved that given exchangeable demands and exclusive use of on-demand staffing in the system-wide optimal solution, $w_i^F \geq 0 \forall i$ as long as the aggregate demand has a log-concave density with no probability mass at zero. Our numerical study shows that $w_i^F \geq 0$ is satisfied in most cases, even for nonexchangeable demands and mixed staffing strategies.

the system when (i) the platform transfers a contingent payment $t_i = r_0 - w_i^C(A_i, \hat{X}_i)$ to a third party, (ii) contract parameters m_i 's are chosen according to (14), and (iii) membership fees r_i 's and $r_0 \geq 0$ satisfy (17) and

$$\sum_{i \in I} r_i - nr_0 \geq C(\mathbf{Q}^*, \eta^*) - p \sum_{i \in I} (1 - m_i) \mathbb{E} [(X_i - Q_i^*)^+]. \quad (20)$$

Furthermore, the third party receives an expected payoff of nr_0 .

Inequality (20) is a simple modification of the platform's participation constraint (17) since the platform operator pays an extra amount of nr_0 to engage the third party. If the third party is a nonprofit organization, one can set $r_0 = 0$ such that all the benefits of coordination are distributed among the platform operator and firms.

For the cases with $w_i^F < 0$ for some $i \in I$, we propose a simple modified fixed-rate mechanism: We use $\max\{w_i^F, 0\}$ as an approximate fixed rate instead of w_i^F while the third party transfer is modified to $t_i = r_0 - w_i^C(A_i, \check{X}_i) - \min\{w_i^F, 0\}$. This modification ensures that each firm will truthfully report as all the fixed rates are kept nonnegative and that the platform's cash flow is unchanged so that its incentive remains aligned with the system optimality. The modification has some impact on Firm i 's staffing decision, as it is equivalent to modifying parameter m_i to $\min\{m_i, \beta_i^*\}$. Nevertheless, the impact appears rather minor. We have resolved our numerical instances with $w_i^F < 0$ under this modified fixed-rate mechanism and observed that the loss of system efficiency due to the above modification is only 0.04% on average and 0.40% in the worst case. Such a minor loss also provides sufficient gain from implementing our mechanism for the system to engage a third party through an extra payment r_0 .

8 Discussion and Conclusion

In this paper, we study the management of an on-demand workforce via a B2B staffing platform, serving cost-minimizing companies. We show that employers can reduce permanent workers and lower staffing costs with access to an on-demand platform. Despite the platform's need to interact with self-interested employers, we prove that the benefits of an on-demand workforce can be maximized at the system-wide optimal level through novel contracts based on fill-rate targets. The proposed contract mechanism shares key features of the performance-based contracting used in temporary staffing services, among many other service industries such as healthcare, public sector services, and aircraft after-sale services (Petersen et al., 2006; Heinrich and Choi, 2007; Guajardo et al., 2012). Our results suggest that to coordinate an on-demand workforce sharing system, the contracts should anchor on the fill rates—a key measure of service level in such

system—which decide the fixed per-unit labor rates. There should also be a contingent payment scheme that rewards exceptional service and penalizes underperformance according to contracted service levels. These are exactly the features of performance based-contracting, which shows that our proposed contract mechanism is practically feasible and should be considered for on-demand staffing services. Moreover, since many on-demand staffing platforms have realized the value of high fill rates for their clients and started to advertise this as a key advantage, it is in the platforms’ interest to explore the proposed fill-rate based contract and allocation mechanisms. Our results provide preliminary analysis and support for these mechanisms in on-demand workforce services while also contributing to the literature on performance-based contracting by considering multiple clients and incorporating a fill-rate-based allocation mechanism.

Throughout the paper, we have assumed that the platform operator can only control the compensation rate. If the platform can influence the size of potential on-demand workers (e.g., by imposing a cap on the number of active workers), the system-wide optimal strategy could be more involved. However, provided the system-wide optimal solution can be evaluated beforehand, the proposed coordination mechanism remains valid because, under the proposed contract, the platform operator’s profit-maximizing problem is aligned with the system-wide cost-minimizing problem (Section 5.2.1), whether the size of potential workers is a decision variable or not. We have made some other assumptions in the model, which suggest potential directions for future research. We have focused on a homogeneous workforce and job positions, suitable for jobs requiring similar skills. In some cases, platforms may offer specialized workers for different job positions. An interesting extension would be to consider multiple types of jobs and a pool of heterogeneous workers; each worker can perform a subset of job types. Then, the staffing system would resemble a network of flexible capacity studied in the literature (e.g., Lyu et al., 2019), making the allocation and contracting problem more challenging. We have also assumed that the on-demand and permanent workforce yields the same level of productivity. However, for tasks requiring significant proficiency or commitment, permanent employees’ productivity may differ from on-demand workers, as empirically studied by Kesavan et al. (2014). This could be modeled by vertically differentiated workers with varying productivity levels. Lastly, we considered risk-neutral firms and assumed cost parameters as common knowledge. Coordinating a system with private cost information or risk-averse firms remains an open question for future research.

Acknowledgments

The authors thank Prof. Guillaume Roels (Department Editor), the Associate Editor and two anonymous reviewers for their constructive feedback. The authors are grateful to Prof. Chung-Piaw Teo for helpful discussions on the earlier version of the paper, and to Girbien de Bruin (Director of Industry and Logistics at Tempo-Team) for sharing her knowledge as to

the practice of performance-based contracting in the staffing industry.

References

- Alptekinoglu, A., A. Banerjee, A. Paul, N. Jain. 2013. Inventory pooling to deliver differentiated service. *Manufacturing & Service Operations Management* **15**(1) 33–44.
- Anupindi, R., Y. Bassok. 1999. Centralization of stocks: Retailers vs. manufacturer. *Management Science* **45**(2) 178–191.
- Asadpour, Arash, Xuan Wang, Jiawei Zhang. 2020. Online resource allocation with limited flexibility. *Management Science* **66**(2) 642–666.
- Bagnoli, M., T. Bergstrom. 2005. Log-concave probability and its applications. *Economic theory* **26**(2) 445–469.
- Benjaafar, Saif, Jian-Ya Ding, Guangwen Kong, Terry Taylor. 2021. Labor welfare in on-demand service platforms. *Manufacturing & Service Operations Management* .
- Benjaafar, Saif, Ming Hu. 2019. Operations management in the age of the sharing economy: What is old and what is new? *Forthcoming, Manufacturing and Service Operations Management* .
- Benjaafar, Saif, Shihong Xiao, Xiaotang Yang. 2020. Do workers and customers benefit from competition between on-demand service platforms? *Available at SSRN* .
- Bhattacharyya, Sugato, Francine Lafontaine. 1995. Double-sided moral hazard and the nature of share contracts. *The RAND Journal of Economics* 761–781.
- Cachon, G. P. 2003. Supply chain coordination with contracts. *Handbooks in operations research and management science* **11** 227–339.
- Cachon, G. P., K. M. Daniels, R. Lobel. 2017. The role of surge pricing on a service platform with self-scheduling capacity. *Manufacturing & Service Operations Management* **19**(3) 368–384.
- Cachon, G. P., M. A. Lariviere. 1999. Capacity choice and allocation: Strategic behavior and supply chain performance. *Management Science* **45**(8) 1091–1108.
- Cachon, G. P., M. A. Lariviere. 2005. Supply chain coordination with revenue-sharing contracts: strengths and limitations. *Management science* **51**(1) 30–44.
- Cachon, Gerard P. 2001. Stock wars: inventory competition in a two-echelon supply chain with multiple retailers. *Operations Research* **49**(5) 658–674.
- Chen, F., A. Federgruen, Y.-S. Zheng. 2001. Coordination mechanisms for a distribution system with one supplier and multiple retailers. *Management Science* **47**(5) 693–708.
- Chen, S., H. Lee, K. Moynadeh. 2016. Supply chain coordination with multiple shipments: the optimal inventory subsidizing contracts. *Operations Research* **64**(6) 1320–1337.
- Chen, Shi, Hau Lee. 2016. Incentive alignment and coordination of project supply chains. *Management Science* **63**(4) 1011–1025.
- Chen, Y, M. Hu. 2019. Pricing and matching in the sharing economy. *Sharing Economy*. Springer, 137–164.
- Corbett, C. J., K. Rajaram. 2006. A generalization of the inventory pooling effect to nonnormal dependent demand. *Manufacturing & Service Operations Management* **8**(4) 351–358.
- Dong, L., N. Rudi. 2004. Who benefits from transshipment? exogenous vs. endogenous wholesale prices. *Management Science* **50**(5) 645–657.
- Eppen, G. D. 1979. Effects of centralization on expected costs in a multi-location newsboy problem. *Management Science* **25**(5) 498–501.
- Gallego, Guillermo, Paul Zipkin. 1999. Stock positioning and performance estimation in serial production-transportation systems. *Manufacturing & Service Operations Management* **1**(1) 77–88.
- Guajardo, J. A., M. A Cohen, S.-H. Kim, S. Netessine. 2012. Impact of performance-based contracting on product reliability: An empirical analysis. *Management Science* **58**(5) 961–979.
- Gurvich, I., M. Lariviere, A. Moreno. 2019. Operations in the on-demand economy: Staffing services with self-scheduling capacity. *Sharing Economy*. Springer, 249–278.
- Hasija, S., E. J. Pinker, R. A. Shumsky. 2008. Call center outsourcing contracts under information asymmetry. *Management Science* **54**(4) 793–807.

- Heinrich, C. J., Y. Choi. 2007. Performance-based contracting in social welfare programs. *The American Review of Public Administration* **37**(4) 409–435.
- Hu, B., I. Duenyas, D. R. Beil. 2013. Does pooling purchases lead to higher profits? *Management Science* **59**(7) 1576–1593.
- InStaff. 2016. Staffing-as-a-service platform for the €170 billion short term staffing market. Available online at: http://www.instaff.jobs/download/InStaff_Pitch_Deck.pdf. Retrieved at 2018-01-26.
- Kesavan, S., B. R. Staats, W. Gilland. 2014. Volume flexibility in services: The costs and benefits of flexible labor resources. *Management Science* **60**(8) 1884–1906.
- Kim, S.-H., M. A. Cohen, S. Netessine. 2007. Performance contracting in after-sales service supply chains. *Management Science* **53**(12) 1843–1858.
- Lin, Xiaogang, Tao Lu, Xin Wang. 2022. Mergers between on-demand service platforms Working paper, available at SSRN: <https://ssrn.com/abstract=3251761>.
- Lyu, Guodong, Wang-Chi Cheung, Mabel C Chou, Chung-Piaw Teo, Zhichao Zheng, Yuanguang Zhong. 2019. Capacity allocation in flexible production networks: Theory and applications. *Management Science* **65**(11) 5091–5109.
- Mak, Ho-Yin, Zuo-Jun Max Shen. 2014. Pooling and dependence of demand and yield in multiple-location inventory systems. *Manufacturing & Service Operations Management* **16**(2) 263–269.
- Netessine, S., N. Rudi. 2006. Supply chain choice on the internet. *Management Science* **52**(6) 844–864.
- Pasternack, B. A. 1985. Optimal pricing and return policies for perishable commodities. *Marketing science* **4**(2) 166–176.
- Petersen, L. A., L. D. Woodard, T. Urech, C. Daw, S. Sookanan. 2006. Does pay-for-performance improve the quality of health care? *Annals of internal medicine* **145**(4) 265–272.
- Ren, Z. J., Y.-P. Zhou. 2008. Call center outsourcing: Coordinating staffing level and service quality. *Management Science* **54**(2) 369–383.
- Robbins, Thomas R, Terry P Harrison. 2011. New project staffing for outsourced call centers with global service level agreements. *Service Science* **3**(1) 41–66.
- Shang, K. H., J.-S. Song, P.H. Zipkin. 2009. Coordination mechanisms in decentralized serial inventory systems with batch ordering. *Management Science* **55**(4) 685–695.
- SIA. 2020. The gig economy and talent platform landscape 2020. Posted 4 September 2020; available online at: <https://www2.staffingindustry.com/Research/Research-Reports/Americas/The-Talent-Platform-Landscape-2020-Update>.
- SIA. 2021a. Nurse staffing: Online platforms poised to make their mark. Posted 13 January 2021; available online at: <https://www2.staffingindustry.com/row/Editorial/Healthcare-Staffing-Report/Jan.-14-2021/Nurse-staffing-Online-platforms-poised-to-make-their-mark>.
- SIA. 2021b. The talent platform landscape: 2021 update. Posted 17 September 2021; available online at: <https://www2.staffingindustry.com/Research/Research-Reports/Americas/The-Talent-Platform-Landscape-2021-Update>.
- Strauss Law Blog. 2015. New york workers subject to illegal “on-demand” staffing may be entitled to compensation. Posted 27 April 2015; available online at: <http://www.strauslawpllc.com/blog/2015/4/27/new-york-workers-subject-to-illegal-on-demand-staffing-may-be-entitled-to-compensation>.
- Swaminathan, J. M., R. Srinivasan. 1999. Managing individual customer service constraints under stochastic demand. *Operations research letters* **24**(3) 115–125.
- Taylor, M. 2017. How important is fill ratio for staffing companies. Available at: <https://www.hrzone.com/community/blogs/mark-taylor/how-important-is-fill-ratio-for-staffing-companies>. Retrieved on 2018-01-28.
- Taylor, T. A. 2018. On-demand service platforms. *Manufacturing & Service Operations Management* **20**(4) 704–720.
- Ülkü, S., L. B. Toktay, E. Yücesan. 2007. Risk ownership in contract manufacturing. *Manufacturing & Service Operations Management* **9**(3) 225–241.
- Zhang, Can, Atalay Atasu, Turgay Ayer, L Beril Toktay. 2020. Truthful mechanisms for medical surplus product allocation. *Manufacturing & Service Operations Management* **22**(4) 735–753.
- Zhang, J. 2003. Managing multi-customer service level requirements with a simple rationing policy. *Operations Research Letters* **31**(6) 477–482.
- Zhong, Y., Z. Zheng, M. C. Chou, C.-P. Teo. 2017. Resource pooling and allocation policies to deliver differentiated service. *Management Science* **64**(4) 1555–1573.
- Zipkin, Paul Herbert. 2000. *Foundations of inventory management*.

Appendices for “Maximizing the Benefits of an On-Demand Workforce: Fill-Rate-Based Allocation and Coordination Mechanisms”

In this document, we present some supplemental content referenced in the main paper. The detailed proofs of all our results are given in Appendix A.

A Technical Proofs

A.1 Proof of Theorem 1

In general, $C(\mathbf{Q}, \eta)$ is not convex in η . We therefore change the decision variable η to S by invoking the relation $S = NG(\eta)$. We define

$$\hat{C}(\mathbf{Q}, S) = \hat{C}\left(\mathbf{Q}, G^{-1}\left(\frac{S}{N}\right)\right) = p \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i)^+ - S \right)^+ \right] + G^{-1}\left(\frac{S}{N}\right) S + c \sum_{i \in I} Q_i,$$

and the problem is transformed to

$$\begin{aligned} \min_{\mathbf{Q}, S} \quad & \hat{C}(\mathbf{Q}, S) \\ \text{s.t.} \quad & Q_i \geq 0, \forall i \in I \\ & 0 \leq \bar{S} \leq N \end{aligned}$$

The following claim is true by the log-concavity of G .

Claim 1. $\hat{C}(\mathbf{Q}, S)$ is jointly convex in \mathbf{Q} and S .

Proof of Claim 1. The term $G^{-1}(S/N)\bar{S}$ is convex in each S . To see this, taking the derivative with respect to S yields

$$\frac{\partial}{\partial S} \left[G^{-1}\left(\frac{S}{N}\right) S \right] = G^{-1}\left(\frac{S}{N}\right) + \frac{S}{N} \frac{1}{g\left(G^{-1}\left(\frac{S}{N}\right)\right)}.$$

Define $\phi(x) := x + G(x)/g(x)$. By the log-concavity of G , $\phi(x)$ is increasing x . The convexity of $G^{-1}(S/N)S$ follows by noticing that $\partial[G^{-1}(S/N)S]/\partial S = \phi(G^{-1}(S/N))$ and $G^{-1}(S/N)$ is an increasing function of S .

It remains to show that the first term of function $\hat{C}(\mathbf{Q}, S)$ is jointly convex, since the last term is linear. Note that $\sum_{i \in I} (X_i - Q_i)^+ - \bar{S}$ is clearly jointly convex in \mathbf{Q} and S . Furthermore, the operator $(x)^+$ is nondecreasing and convex. It is well known that their composition is still convex. As convexity is preserved under expectation and summation, the first term of function $\hat{C}(\mathbf{Q}, S)$ is indeed jointly convex. \square

From Lemma 1, the first-order conditions are necessary and sufficient for optimality. It then suffices to verify that in the optimal solution $Q_i = Q_i^*$ and $S = S^* = NG(\eta^*)$ as described in Theorem 1, the first-order conditions hold, i.e., $\partial\hat{C}/\partial S = 0$ (or $\partial\hat{C}/\partial S \geq 0$, if $S^* = 0$) and $\partial\hat{C}/\partial Q_i = 0$ (or $\partial\hat{C}/\partial Q_i \geq 0$, if $Q_i^* = 0$).

For part (i), we verify that $S = 0$ and $\mathbf{Q} = \mathbf{Q}^p$ satisfy:

$$\begin{aligned} \frac{\partial \hat{C}}{\partial S} \Big|_{S=0, \mathbf{Q}=\mathbf{Q}^p} &= -p\mathbb{P}\left(0 < \sum_{i \in I} (X_i - Q_i^p)^+\right) + \underline{\tau} \\ &= -p\mathbb{P}(X_i > Q_i^p, \text{ for some } i \in I) + \underline{\tau} \geq 0, \end{aligned} \quad (\text{A.1})$$

and

$$\begin{aligned} \frac{\partial \hat{C}}{\partial Q_i} \Big|_{S=0, \mathbf{Q}=\mathbf{Q}^p} &= -p\mathbb{P}\left(Q_i^p < X_i, 0 < \sum_{i \in I} (X_i - Q_i^p)^+\right) + c \\ &= -p\mathbb{P}(Q_i^p < X_i) + c = 0. \end{aligned} \quad (\text{A.2})$$

Inequality (A.1) holds if and only if the condition in part (ii) is satisfied: $\underline{\tau} \geq \phi^p = p\{1 - \mathbb{P}(X_i \leq Q_i^p, \forall i)\}$, and equality (A.2) holds by the definition of Q_i^p .

For part (ii), by the optimality condition, we need to verify that $\mathbf{Q} = 0$ and $S = NG(\eta^\circ)$ are such that

$$\frac{\partial \hat{C}}{\partial S} \Big|_{S=NG(\eta^\circ), \mathbf{Q}=0} = -p\mathbb{P}\left(NG(\eta^\circ) < \sum_{i \in I} X_i\right) + \phi(\eta^\circ) = 0, \quad (\text{A.3})$$

where $\phi(x) := x + G(x)/g(x)$, and

$$\frac{\partial \hat{C}}{\partial Q_i} \Big|_{S=NG(\eta^\circ), \mathbf{Q}=0} = -p\mathbb{P}\left(NG(\eta^\circ) < \sum_{i \in I} X_i\right) + c \geq 0. \quad (\text{A.4})$$

Equality (A.3) holds by the definition of η° . Moreover, it implies that $p\mathbb{P}(NG(\eta^\circ) < \sum_{i \in I} X_i) = \phi(\eta^\circ) = \eta^\circ + G(\eta^\circ)/g(\eta^\circ)$. Substituting this equality into $\partial \hat{C}/\partial Q_i$, one can see that inequality (A.4) holds if and only if $c \geq \phi^\circ$, which completes the proof of part (ii).

Note that the proofs of parts (i) and (ii) provide the sufficient and also necessary conditions for the optimality of boundary solutions $S^* = 0$ and $\mathbf{Q}^* = 0$, respective. Therefore, if $c < \phi^\circ$ or $\underline{\tau} \geq \phi^p$, the interior solution must be attained and the equations in part (iii) follow immediately by equating the first-order derivatives of \hat{C} with respect to S and the Q_i 's to zero.

A.2 Proof of Corollary 1

When X_i 's are independent across i , we have

$$\phi^p = p \left(1 - \prod_{i \in I} F_i(Q_i^p)\right) = p \left[1 - \left(\frac{p-c}{p}\right)^n\right].$$

The statements in Part (i) immediately follow. By Theorem 1, the on-demand workforce will always be used in the optimal solution as long as $\underline{\tau} < p$ —i.e., it is profitable by itself, as we have postulated.

Note that $\mathbf{X} \leq_{SM} \mathbf{X}'$ implies that \mathbf{X} and \mathbf{X}' have the same univariate marginals. Hence, the values for Q_i^p 's are the same for two demand vectors (Shaked and Shanthikumar, 1997). By definition, $\mathbf{X} \leq_{SM} \mathbf{X}'$ if and only if $\mathbb{E}[f(\mathbf{X})] \leq \mathbb{E}[f(\mathbf{X}')] for all supermodular functions f such that the expectations exist. Note that $1 - \mathbb{P}\{X_i \leq Q_i^p, \forall i\} = 1 - \mathbb{E}[\prod_{i \in I} \mathbb{I}(X_i \leq Q_i^p)]$, where $\mathbb{I}(\cdot)$ represents the indicator$

function. Thus, it suffices to show that the function $f(x_1, x_2, \dots, x_n) = \prod_{i \in I} \mathbb{I}(x_i \leq Q_i^p)$ is supermodular in (x_1, x_2, \dots, x_n) .

To see this, one can verify that the following inequality holds for any two vectors \mathbf{x} and \mathbf{x}' ,

$$f(\mathbf{x}) + f(\mathbf{x}') \leq f(\mathbf{x} \wedge \mathbf{x}') + f(\mathbf{x} \vee \mathbf{x}'),$$

where

$$\mathbf{x} \wedge \mathbf{x}' = (\min\{x_1, x'_1\}, \min\{x_2, x'_2\}, \dots, \min\{x_n, x'_n\}), \text{ and}$$

$$\mathbf{x} \vee \mathbf{x}' = (\max\{x_1, x'_1\}, \max\{x_2, x'_2\}, \dots, \max\{x_n, x'_n\}).$$

If $f(\mathbf{x}) = f(\mathbf{x}') = 1$, then $x_i \leq Q_i^p$ and $x'_i \leq Q_i^p$ for all $i \in I$. Thus, $f(\mathbf{x} \wedge \mathbf{x}') = f(\mathbf{x} \vee \mathbf{x}') = 1$. If $f(\mathbf{x}) = 1$ and $f(\mathbf{x}') = 0$, then $x_i \leq Q_i^p$ for all $i \in I$, but there exists some j such that $x'_j > Q_j^p$. As a result, $f(\mathbf{x} \wedge \mathbf{x}') = 1$ and $f(\mathbf{x} \vee \mathbf{x}') = 0$. The inequality holds with equality for the above two cases. Finally, if $f(\mathbf{x}) = f(\mathbf{x}') = 0$, the inequality always hold and can sometimes be strict. For instance, when $n = 2$, $x_1 \leq Q_1^p$, $x_2 > Q_2^p$, $x'_1 > Q_1^p$ and $x_2 \leq Q_2^p$, we have $f(\mathbf{x} \wedge \mathbf{x}') = 1$ and $f(\mathbf{x} \vee \mathbf{x}') = 0$.

A.3 Proof of Corollary 2

Note that $\hat{C}(\mathbf{Q}, S)$ is supermodular in S and Q_i (for any i), or equivalently, $\partial^2 C / (\partial S \partial Q_i) \geq 0$. In a convex minimization problem, it follows that the optimal choice of Q_i is nonincreasing in S , and therefore $Q_i^p \geq \hat{Q}_i^*$.

A.4 Proof of Lemma 1

We prove this lemma in three steps. In the first step, we derive a property of the allocation from the TFRB policy when the target fill rates satisfy the constraints of problem (8) under any given workforce level S and firms' permanent staffing levels \mathbf{Q} . Next, we show that the system-wide minimum staffing cost can be achieved by the TFRB policy when the workforce level and firms' permanent staffing levels are set to the optimal ones in the centralized system (as in Theorem 1). Lastly, we strengthen the property shown in the first step for optimal target fill rates from solving problem (8).

A.4.1 Step 1 in Proof of Lemma 1

We first prove that the expected allocation satisfies $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)] \geq \beta_i \mathbb{E}[(X_i - Q_i)^+]$ for all $i \in I$ for any feasible β of problem (8) under any given S and \mathbf{Q} . The proof in this step follows from Theorem 3 in Zhong et al. (2017). For completeness, we adapt the proof to our context and present it below.

To this end, we first show that the TFRB policy is able to deliver the expected fill rates in the multi-period setting as described in Algorithm 1. Specifically, we will show that the TFRB policy satisfies the sufficient condition of Blackwell's approachability theorem (Blackwell, 1956) such that the target fill rates are approachable in the multi-period setting.

At the beginning of simulated period t , the TFRB policy uses the gap between the expected allocation in the past $(t - 1)$ simulated periods and the actual allocation received, denoted as $\mathbf{A}(t) =$

$(A_1(t), A_2(t), \dots, A_n(t))$, to determine the allocation sequence in simulated period t . Following the terminology of Blackwell's approachability theorem, the *reward* gained by Firm i in simulated period t , $R_i(t)$, is defined as $\beta_i \mathbb{E}[\hat{X}_i] - A_i(t)$. The reward of Firm i is negative when the current-period allocation exceeds the target fill rate β_i , and it is positive when the target fill rate is not satisfied. The *debt* of Firm i at the beginning of simulated period t before allocation, $r_i(t-1)$, is the gap described in Algorithm 1 to prioritize the allocation in each period, i.e.,

$$r_i(t-1) := (t-1)\beta_i \mathbb{E}[\hat{X}_i(t)] - \sum_{k=1}^{t-1} A_i(k).$$

The time-average debt from period 1 up to period $(t-1)$ is $\rho_i(t-1) = r_i(t-1)/(t-1)$. After the allocation in period t , the time-average debt of Firm i becomes

$$\rho_i(t) = \frac{r_i(t)}{t} = \frac{r_i(t-1) + R_i(t)}{t}.$$

Let \mathcal{D} denote the set of nonpositive orthant in \mathbb{R}^n . According to Blackwell's approachability theorem, we want the time-average debt to approach $\mathcal{D} := \{z = [z_1, z_2, \dots, z_n] \mid z_i \leq 0, \forall i \in I\}$ to ensure that Firm i receives an expected fill rate of at least β_i .

Suppose that at the beginning of some period t , the time-average debt, $\boldsymbol{\rho}(t-1) = (\rho_1(t-1), \rho_2(t-1), \dots, \rho_n(t-1))$ is not in \mathcal{D} . For ease of notation, let $\boldsymbol{\alpha} = \boldsymbol{\rho}(t-1)$. Without loss of generality, we assume that $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_m > 0 \geq \alpha_{m+1} \geq \alpha_{m+2} \geq \dots \geq \alpha_n$. The point in \mathcal{D} closest to $\boldsymbol{\alpha}$ is $\boldsymbol{\gamma} = (0, 0, \dots, 0, \alpha_{m+1}, \alpha_{m+2}, \dots, \alpha_n)$. The hyperplane perpendicular to the line segment $\boldsymbol{\alpha}\boldsymbol{\gamma}$ is $H := \{z \in \mathbb{R}^n : \sum_{i=1}^m z_i \alpha_i = 0\}$. Note that $\boldsymbol{\alpha} \in H^+ := \{z \in \mathbb{R}^n : \sum_{i=1}^m z_i \alpha_i > 0\}$. We only need show that when the allocation follows the TFRB policy, the mean reward $(\mathbb{E}[R_1(t)], \mathbb{E}[R_2(t)], \dots, \mathbb{E}[R_n(t)])$ lies in $H^- \cup H := \{z \in \mathbb{R}^n : \sum_{i=1}^m z_i \alpha_i \leq 0\}$, where $H^- := \{z \in \mathbb{R}^n : \sum_{i=1}^m z_i \alpha_i < 0\}$. We make the following claim and continue the proof first.

Claim 2.

$$\varpi_j := \sum_{i=1}^j \mathbb{E}[R_i(t)] = \mathbb{E} \left[\sum_{i=1}^j \beta_i \mathbb{E}[\hat{X}_i] - \sum_{i=1}^j A_i(t) \right] \leq 0, \forall j = 1, 2, \dots, m.$$

From the above claim, we have $\varpi_j \alpha_j \leq 0$, for all $j = 1, 2, \dots, m$. Therefore,

$$\begin{aligned} \sum_{i=1}^m \mathbb{E}[R_i(t)] \alpha_i &= \varpi_1 \alpha_1 + \sum_{i=2}^m (\varpi_i - \varpi_{i-1}) \alpha_i \\ &\leq \varpi_1 \alpha_1 + \sum_{i=2}^m \varpi_i \alpha_i - \sum_{i=2}^m \varpi_{i-1} \alpha_{i-1} \\ &= \varpi_m \alpha_m \\ &\leq 0, \end{aligned}$$

where the second inequality follows from $\alpha_{i-1} \geq \alpha_i$ and Claim 2. By Blackwell's approachability theorem, our allocation policy makes the time-average debt to approach \mathcal{D} . Hence, the long-run average performance of the TFRB policy attains the desired expected fill-rate requirements for all firms, i.e.,

$\liminf_{T \rightarrow \infty} \sum_{t=1}^T A_i(t)/T \geq \beta_i \mathbb{E}[\hat{X}_i]$ almost surely for all $i \in I$.

Following Step 1 of the TFRB policy, the priority list L is randomly drawn from the priority lists used in the multi-period setting, i.e., $L(2), L(3), \dots, L(T)$. Then the expected amount received by Firm i following the priority list L in our original single-period problem satisfies

$$\mathbb{E} \left[A_i^L(\hat{\mathbf{X}}, S) \right] = \mathbb{E} \left[\frac{\sum_{t=2}^T A_i^{L(t)}(\hat{\mathbf{X}}, S)}{T-1} \right] = \mathbb{E} \left[\frac{\sum_{t=2}^T A_i^{L(t)}(t)}{T-1} \right] \geq \beta_i \mathbb{E}[\hat{X}_i], \text{ for large enough } T,$$

where the superscripts are used to highlight the specific priority lists used in the allocation; the first equality follows from the random draw of the priority list; the second equality holds because the multi-period stochastic demands are independent and identically distributed and follow the same distribution as $\hat{\mathbf{X}}$; and the last inequality is the approachability result just shown above for the multi-period setting. Note that the allocation for the first period does not affect the limiting behavior of the average allocation as T goes to infinity.

Proof of Claim 2. Now we are left to prove the claim. There are two possible cases when the allocation process stops, which we discuss separately. In the first case, the allocation stops at the $(k+1)$ th firm in the priority list. According to the TFRB policy, this will happen only when the remaining capacity, $S - \sum_{i=1}^k A_i(t)$, is less than 0, i.e., $A_{k+1}(t) = 0$ if $\sum_{i=1}^{k-1} A_i(t) \leq S$ and $\sum_{i=1}^k A_i(t) \geq S$. For any $i \leq \min\{k, m\}$, $A_i(t) = \hat{X}_i(t)$. Then we have

$$\mathbb{E} \left[\sum_{i=1}^j A_i(t) \right] = \sum_{i=1}^j \mathbb{E}[\hat{X}_i] \geq \sum_{i=1}^j \beta_i \mathbb{E}[\hat{X}_i], \forall j = 1, 2, \dots, \min\{k, m\}.$$

If $m \leq k$, the claim is proved. If $m > k$,

$$\sum_{i=1}^j A_i(t) \geq S \geq \min \left\{ S, \sum_{i=1}^j \hat{X}_i(t) \right\}, \forall j = k+1, k+2, \dots, m.$$

Taking the expectations on both sides of the above inequality and applying the constraints in problem (8), we have

$$\mathbb{E} \left[\sum_{i=1}^j A_i(t) \right] \geq \mathbb{E} \left[\min \left\{ S, \sum_{i=1}^j \hat{X}_i(t) \right\} \right] \geq \sum_{i=1}^j \beta_i \mathbb{E}[\hat{X}_i], \forall j = k+1, k+2, \dots, m,$$

which implies the inequalities stated in the claim.

In the second case, the allocation process does not stop after going through all of the firms on the priority list once. According to the TFRB policy, it is obvious that for any Firm i , $A_i(t) \geq \beta_i \hat{X}_i(t)$. The claim follows immediately after taking the expectations. \square

A.4.2 Step 2 in Proof of Lemma 1

To demonstrate that the system-wide optimality is achieved when the TFRB policy is implemented and when \mathbf{Q} and η are set as the system-wide optimal solution, refer to Section 4, where the objective is to

minimize the total cost expressed as:

$$p \sum_{i \in I} \mathbb{E} \left[(X_i - Q_i)^+ - A_i \left((\mathbf{X} - \mathbf{Q})^+, S \right) \right] + \eta NG(\eta) + c \sum_{i \in I} Q_i.$$

According to the TFRB policy, the allocation process stops either when the workforce S is depleted or when all demands are met. Consequently, the no-waste condition in (4) is fulfilled, i.e., $\sum_{i \in I} A_i((\mathbf{X} - \mathbf{Q})^+, S) = \min\{S, \sum_{i \in I} (X_i - Q_i)^+\}$. As explained in Section 4, this equation causes the allocation terms, $A_i(\cdot)$'s, to cancel out in the total cost function. Therefore, the system-wide staffing cost is minimized when $\mathbf{Q} = \mathbf{Q}^*$ and $\eta = \eta^*$.

A.4.3 Step 3 in Proof of Lemma 1

Lastly, we prove the ‘‘moreover’’ part of the lemma. That is, the inequalities shown in Step 1 are actually tight for the optimal β^* of problem (8), i.e., $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)] = \beta_i^* \mathbb{E}[(X_i - Q_i)^+]$ for all $i \in I$ under any given S and \mathbf{Q} , which is a stronger version of the last part of Lemma 1.

To deal with the set of constraints in problem (8), we introduce several useful concepts. Given a finite set Ω , define the power set of Ω as $2^\Omega = \{U : U \subseteq \Omega\}$. A set function $g : 2^\Omega \rightarrow \mathbb{R}$ is *submodular* if $g(A \cup B) + g(A \cap B) \leq g(A) + g(B)$ for all $A, B \subseteq \Omega$. The function $g : 2^\Omega \rightarrow \mathbb{R}$ is called a *rank function* if it satisfies (1) $g(\emptyset) = 0$, (2) $g(A) \leq g(B)$ whenever $A \subseteq B \subseteq \Omega$, and (3) g is submodular. The polyhedron,

$$P(g, \Omega) = \left\{ x \in \mathbb{R}_+^{|\Omega|} : \sum_{i \in U} x_i \leq g(U), \forall U \subseteq \Omega \right\},$$

is called a *polymatroid* if g is a rank function, where $|\Omega|$ denotes the cardinality of Ω . We recall the following result from Edmonds (1970), which states that when optimizing a linear function over a polymatroid, an optimal solution can be found using a greedy method.

Claim 3 (Edmonds (1970)). Let $n = |\Omega|$. For any vector $a \in \mathbb{R}^n$, where $a_{j(1)} \geq a_{j(2)} \geq \dots \geq a_{j(k)} > 0 \geq a_{j(k+1)} \dots \geq a_{j(n)}$ for some permutation $\{j(1), j(2), \dots, j(n)\}$ of Ω . $\sum_{i \in I} a_i x_i$ is maximized over $x \in P(g, \Omega)$ by the vector x^* defined as follows:

$$\begin{aligned} x_{j(1)}^* &= g(\{j(1)\}), \\ x_{j(i)}^* &= g(\{j(1), j(2), \dots, j(i)\}) - g(\{j(1), j(2), \dots, j(i-1)\}), 2 \leq i \leq k, \\ x_{j(i)}^* &= 0, \text{ for } k+1 \leq i \leq n. \end{aligned}$$

Define a set function $g(U) = \mathbb{E}[\min\{S, \sum_{i \in U} (X_i - Q_i)^+\}]$. It is straightforward to verify that: (1) $g(\emptyset) = 0$; (2) $g(U)$ is nondecreasing in U ; and (3) $g(U)$ is submodular. Since $\mathbb{E}[(X_i - Q_i)^+]$'s are nonnegative constants given Q_i 's, problem (8) is equivalent to maximize over $\hat{\beta}_i$'s, where $\hat{\beta}_i = \beta_i \mathbb{E}[(X_i - Q_i)^+]$; that is, we can replace the original decision variables β_i 's with $\hat{\beta}_i$'s. By Claim 3, the optimal $\hat{\beta}_i^*$

can be found in a greedy manner. Consequently, we have

$$\sum_{i \in I} \beta_i^* \mathbb{E}[(X_i - Q_i)^+] = \mathbb{E} \left[\min \left\{ S, \sum_{i \in I} (X_i - Q_i)^+ \right\} \right]. \quad (\text{A.5})$$

Recall that $\sum_{i \in I} A_i((\mathbf{X} - \mathbf{Q})^+, S) = \min\{S, \sum_{i \in I} (X_i - Q_i)^+\}$ on any sample path under the TFRB policy. Therefore, (A.5) implies $\sum_{i \in I} \beta_i^* \mathbb{E}[(X_i - Q_i)^+] = \sum_{i \in I} \mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)]$. Together with $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)] \geq \beta_i \mathbb{E}[(X_i - Q_i)^+]$ for each i from the first step, we deduce that $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)] = \beta_i^* \mathbb{E}[(X_i - Q_i)^+]$ for all $i \in I$.

A.4.4 Remark: Refined Fill Rates by Algorithm 2

We note that our results remain unaffected if we use Algorithm 2 in Section 5.1 to refine the optimal fill rates.

Step 1 of the proof above establishes that the TFRB policy can achieve any feasible fill rates from problem 8 under given S and \mathbf{Q} . Algorithm 2 ensures that the refined fill rates β_i^* 's are feasible under S^* and \mathbf{Q}^* , and thus the TFRB policy guarantees that the target fill rates are met, i.e., $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S^*)] \geq \beta_i^* \mathbb{E}[(X_i - Q_i^*)^+]$ for all i . Step 2 of the proof remains unchanged, as the TFRB policy satisfies the no-waste condition, irrespective of whether the refined fill rates are used or not. Consequently, the system-wide optimality can be achieved as before. Step 3 of the above proof merely demonstrates that the fill rate guarantee $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)] \geq \beta_i^* \mathbb{E}[(X_i - Q_i)^+]$ holds with equality. Recall that the refined fill rates are the solution to a linear maximization over a polymatroid with side constraints $\beta_i \geq b^*$ as studied in Lu and Yao (2008). The grouping algorithm in Lu and Yao (2008) implies that, in the optimal solution, the polymatroid constraint for the entire set holds with equality, i.e., $\sum_{i \in I} \beta_i^* \mathbb{E}[(X_i - Q_i)^+] = \sum_{i \in I} \mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)]$. Since the TFRB policy ensures $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S^*)] \geq \beta_i^* \mathbb{E}[(X_i - Q_i^*)^+]$ for all i as shown in Step 1 of the proof, it follows that $\mathbb{E}[A_i((\mathbf{X} - \mathbf{Q})^+, S)] = \beta_i^* \mathbb{E}[(X_i - Q_i)^+]$ for all i .

A.5 Proof of Lemma 2

The proof follows the discussion in Section 5.2.1 and Lemma 1.

A.6 Proof of Lemma 3

The proof follows the discussion in Section 5.2.2. The characterization of m_i 's is obtained by equating the first-order derivative of $h_i(Q_i)$ at $Q_i = Q_i^*$ to zero.

A.7 Proof of Theorem 2

The proof follows by invoking the notion of Nash equilibrium and Lemmas 2 and 3, as well as the discussion of participation issues in Section 5.2.3.

A.8 Proof of Corollary 3

Given $\mathbf{Q}^* = 0$, the platform operator's profit is given by

$$\Pi = -C(0, \eta^o) + p \sum_{i \in I} (1 - m_i) \mathbb{E}[X_i] = -C(0, \eta^o) + c \sum_{i \in I} \mathbb{E}[X_i],$$

where m_i is chosen according to (14). If $c \sum_{i \in I} \mathbb{E}[X_i] \geq C(0, \eta^o)$, the platform operator has collected a nonnegative profit without membership fees.

When $\mathbf{Q} = 0$, the realized cost function is supermodular in \mathbf{X} by Lemma 2.6.2 of (Topkis, 1998). Therefore, for any two demand vectors $\mathbf{X}' \geq_{sm} \mathbf{X}$, let $\eta^{o'}$ and η^o denote the system-wide optimal compensation rates corresponding to \mathbf{X}' and \mathbf{X} , respectively. It follows that $C(0, \eta^{o'}; \mathbf{X}') \geq C(0, \eta^{o'}; \mathbf{X}) \geq C(0, \eta^o; \mathbf{X})$, where the first inequality is due to the definition of supermodular order and the second inequality holds by the optimality of η^o .

A.9 Proof of Proposition 1

By Theorem 1, when the X_i 's are exchangeable, the system-wide optimal staffing level of each firm is identical, i.e., $Q_i^* = Q^*$ for all $i \in I$. By the definitions of w_i^F and w_i^C , i.e., (11)–(12), and the range of membership fees r_i 's, i.e., (16)–(17), the parameters that make the contract firm-dependent are the m_i 's and β_i^* 's. It suffices to show that m_i and β_i^* can be firm-independent.

Per Lemma 3, the value of m_i reduces to a firm-independent value, $m = (p\bar{F}(Q^*) - c)/(p\bar{F}(Q^*)) = 1 - c/(p\bar{F}(Q^*))$, since the marginal distributions are identical.

It remains to show that we can find a set of identical fill rates, $\beta_i^* = \beta^*$ for all $i \in I$, which is system-wide optimal. By Lemma 1, it is equivalent to proving that problem (8) has an optimal solution such that $\beta^* = (\beta^*, \beta^*, \dots, \beta^*)$.

Lemma 1 implies that the optimal objective value of problem (8) is equal to $\sum_{i \in I} \beta_i^* \mathbb{E}[(X_i - Q^*)^+] = \sum_{i \in I} \mathbb{E}[A_i^*((\mathbf{X} - \mathbf{Q}^*)^+, S^*)] = \mathbb{E}[\min\{S^*, \sum_{i \in I} (X_i - Q^*)^+\}]$, where the last equality follows from the no-waste condition (4).

Now, we restrict the decision variables to $\beta_i = \beta$ for all $i \in I$ and solve the restricted version of problem (8). If the firm-independent optimal solution β^* leads to the same objective value $\mathbb{E}[\min\{S^*, \sum_{i \in I} (X_i - Q^*)^+\}]$, then β^* is optimal to the original LP. Since X_i 's have an identical marginal distribution, we can write $H := \mathbb{E}[(X_i - Q^*)^+] = \int_{Q^*}^{\infty} (x - Q^*) dF(x)$. Then, the restricted LP can be written as follows:

$$\begin{aligned} \max_{\beta} \quad & n\beta H \\ \text{s.t.} \quad & |U|\beta H \leq \mathbb{E} \left[\min \left\{ S^*, \sum_{i \in U} (X_i - Q^*)^+ \right\} \right], \forall U \subseteq I \\ & \beta \geq 0, \forall i \in I \end{aligned} \tag{A.6}$$

Clearly, the optimal solution is given by

$$\beta^* = \frac{1}{H} \min_{U \subseteq I} \frac{\mathbb{E}[\min\{S^*, \sum_{i \in U} (X_i - Q^*)^+\}]}{|U|},$$

where $|U|$ represents the cardinality of set U . In what follows, we prove that $\mathbb{E}[\min\{S^*, \sum_{i \in U}(X_i - Q^*)^+\}]/|U|$ is nonincreasing in U , thus being minimized when $U = I$.

Consider two sets $U = \{1, 2, \dots, k\}$ and $U' = \{1, 2, \dots, k+1\}$ for any given $k \leq n-1$. We need to evaluate the following difference:

$$\begin{aligned} & \frac{\mathbb{E}[\min\{S^*, \sum_{i \in U'}(X_i - Q^*)^+\}]}{|U'|} - \frac{\mathbb{E}[\min\{S^*, \sum_{i \in U}(X_i - Q^*)^+\}]}{|U|} \\ &= \frac{\mathbb{E}\left[k \min\left\{S^*, \sum_{i=1}^{k+1}(X_i - Q^*)^+\right\} - (k+1) \min\left\{S^*, \sum_{i=1}^k(X_i - Q^*)^+\right\}\right]}{k(k+1)}. \end{aligned}$$

In total, three possible events, denoted by e_1, e_2 and e_3 , could happen for the numerator above.

- (i) Conditional on event $e_1 = \{S^* \leq \sum_{i=1}^k(X_i - Q^*)^+\}$, the numerator $= -S^* \leq 0$.
- (ii) Conditional on event $e_2 = \{\sum_{i=1}^k(X_i - Q^*)^+ < S^* \leq \sum_{i=1}^{k+1}(X_i - Q^*)^+\}$, the numerator $= \mathbb{E}[kS^* - (k+1)\sum_{i=1}^k(X_i - Q^*)^+ | e_2] \leq \mathbb{E}[k\sum_{i=1}^{k+1}(X_i - Q^*)^+ - (k+1)\sum_{i=1}^k(X_i - Q^*)^+ | e_2] = (k(k+1) - k(k+1))\mathbb{E}[(X_i - Q^*)^+ | e_2] = 0$. Note that the last equality holds because X_i 's are exchangeable.
- (iii) Conditional on event $e_3 = \{\sum_{i=1}^{k+1}(X_i - Q^*)^+ < S^*\}$, the numerator $= \mathbb{E}[k\sum_{i=1}^{k+1}(X_i - Q^*)^+ - (k+1)\sum_{i=1}^k(X_i - Q^*)^+ | e_3] = (k(k+1) - (k+1)k)\mathbb{E}[(X_i - Q^*)^+ | e_3] = 0$.

In summary, as the numerator is nonpositive under all the possible events, its total expectation is nonpositive. Hence, $\mathbb{E}[\min\{S^*, \sum_{i \in U'}(X_i - Q^*)^+\}]/|U'| - \mathbb{E}[\min\{S^*, \sum_{i \in U}(X_i - Q^*)^+\}]/|U| \leq 0$. Note that we can arbitrarily re-number the firms, so the above argument applies to any sets $U, U' \subseteq I$ such that $U' = U \cup \{j\}$ where $j \notin U$. Therefore, $\mathbb{E}[\min\{S^*, \sum_{i \in U}(X_i - Q^*)^+\}]/|U|$ is minimized when $U = I$. It follows that the optimal solution $\beta^* = \mathbb{E}[\min\{S^*, \sum_{i \in I}(X_i - Q^*)^+\}]/(nH)$. Substituting β^* into the objective function of the restricted problem (A.6), we have the optimal objective value equal to $\mathbb{E}[\min\{S^*, \sum_{i \in I}(X_i - Q^*)^+\}]$, which is the same as that of the original LP without any restriction on $(\beta_1, \beta_2, \dots, \beta_n)$. Therefore, the firm-independent fill rate $\beta = (\beta^*, \beta^*, \dots, \beta^*)$ also solves the original LP and is thus system-wide optimal.

Invoking $m_i = m$ and $\beta_i^* = \beta^*$ for all $i \in I$ leads to the closed-form expressions of the contract terms.

A.10 Proof of Lemma 4

Similar to the coordination contract discussed in Section 5.2, summing up the contract terms, we have the total payment from Firm i to the platform as $r + p(A_i - \bar{m}\hat{X}_i)$, for any allocated amount A_i and

request $\hat{X}_i = (X_i - Q_i)^+$. For the platform operator, we can transform its expected profit as follows:

$$\begin{aligned}
\Pi(\eta, \mathbf{A}(\cdot)|\mathbf{Q}) &= -\eta NG(\eta) + nr + p \sum_{i \in I} \mathbb{E}[A_i] - p\bar{m} \sum_{i \in I} \mathbb{E}[\hat{X}_i] \\
&= -\eta NG(\eta) + nr + p \mathbb{E} \left[\min \left\{ S, \sum_{i \in I} \hat{X}_i \right\} \right] - p\bar{m} \sum_{i \in I} \mathbb{E}[\hat{X}_i] \\
&= -\eta NG(\eta) - c \sum_{i \in I} Q_i - p \mathbb{E} \left[\left(\sum_{i \in I} \hat{X}_i - S \right)^+ \right] + c \sum_{i \in I} Q_i + p \mathbb{E} \left[\left(\sum_{i \in I} \hat{X}_i - S \right)^+ \right] \\
&\quad + nr + p \mathbb{E} \left[\min \left\{ S, \sum_{i \in I} \hat{X}_i \right\} \right] - p\bar{m} \sum_{i \in I} \mathbb{E}[\hat{X}_i] \\
&= -C(\mathbf{Q}, \eta) + p(1 - \bar{m}) \sum_{i \in I} \mathbb{E}[(X_i - Q_i)^+] + c \sum_{i \in I} Q_i + nr.
\end{aligned}$$

The second equality above follows from the no-waste property of allocation policies; in the third equality, we rearrange the terms to obtain the system cost $C(\mathbf{Q}, \eta) = \eta NG(\eta) + c \sum_{i \in I} Q_i + p \mathbb{E}[(\sum_{i \in I} \hat{X}_i - S)^+]$ in the expression; the fourth equality follows because $(\sum_{i \in I} \hat{X}_i - S)^+ + \min\{S, \sum_{i \in I} \hat{X}_i\} = \sum_{i \in I} \hat{X}_i$. From the above expression of $\Pi(\eta, \mathbf{A}(\cdot)|\mathbf{Q})$, it can be seen that given any \mathbf{Q} , the platform operator will choose an η that minimizes the system cost while being indifferent among all the no-waste allocation policies.

Firm i 's cost function can be written as

$$\begin{aligned}
h_i(Q_i|\mathbf{Q}_{-i}, \mathbf{A}(\cdot), \eta) &= p \mathbb{E} \left[\left((X_i - Q_i)^+ - A_i \right)^+ \right] + p \mathbb{E}[A_i] - p\bar{m} \mathbb{E}[(X_i - Q_i)^+] + r + cQ_i \\
&= p(1 - \bar{m}) \mathbb{E}[(X_i - Q_i)^+] + cQ_i + r.
\end{aligned}$$

Hence, Firm i 's optimal permanent staffing level, denoted by Q_i^I , is the one that minimizes $p(1 - \bar{m}) \mathbb{E}[(X_i - Q_i)^+] + cQ_i + r$. Similar to the case with the coordination contract in Section 5.2, Q_i^I is a dominant strategy of Firm i regardless of the other players' actions (including the platform's allocation policy).

Combining the platform operator's and every firm's best response strategies, we can conclude that in any Nash equilibrium, each Firm i will choose Q_i^I , whereas the platform operator will set the compensation rate as $\eta^I = \arg \min_{\eta} C(\mathbf{Q}^I, \eta)$ while being indifferent among all the no-waste allocation policies.

Let $S^I = NG(\eta^I)$ denote the on-demand staffing level induced by η^I . Similar to the case with our coordination contract, the platform operator will voluntarily adopt the proposed TFRB policy in equilibrium. However, the target fill rates used in the TFRB policy, denoted by $\beta^I = (\beta_1^I, \beta_2^I, \dots, \beta_n^I)$, need to be calculated based on the equilibrium staffing levels \mathbf{Q}^I and $S^I = NG(\eta^I)$, as these staffing levels are different from the system-wide optimal solution. Recall that in the proof of Lemma 1, Steps 1 and 2 are valid for any given \mathbf{Q} and S and thus also applicable to \mathbf{Q}^I and S^I . Hence, β^I can be found

by solving the following linear program:

$$\begin{aligned}
& \max_{\beta} \quad \sum_{i \in I} \beta_i \mathbb{E} [(X_i - Q_i^I)^+] \\
& \text{s.t.} \quad \sum_{i \in U} \beta_i \mathbb{E} [(X_i - Q_i^I)^+] \leq \mathbb{E} \left[\min \left\{ S^I, \sum_{i \in U} (X_i - Q_i^I)^+ \right\} \right], \forall U \subseteq I \\
& \quad \beta_i \geq 0, \forall i \in I
\end{aligned} \tag{A.7}$$

Next, we can use Algorithm 2 to minimize the differences among the β_i^I 's. Note that the parameter b^* in Algorithm 2 serves as a nominal fill rate used only for the calculation of firm-independent contract terms. It differs from the β_i^I 's, the target fill rates used in equilibrium. Since every firm is under the same contract terms in firm-independent contracts, the platform can randomize firm indexes such that each firm will be assigned with any target fill rate β_i^I with equal probability. Then, $\bar{\beta}^I = \sum_{i \in I} \beta_i^I / n$ will be a uniform fill rate expected to be delivered to all firms.

A.11 Proof of Proposition 2

The proposition follows by noticing that $m_i = (p - c)/p$ for all $i \in I$ and that using b^* instead of firm-dependent fill rates β_i^* 's does not affect the platform's staffing and allocation decisions.

A.12 Proof of Lemma 5

Let $A_i(\check{X}_i)$ denote the amount of workforce allocated to Firm i when he reports \check{X}_i , and $h_i(\check{X}_i, \hat{X}_i)$ represent the staffing cost incurred if Firm i reports \check{X}_i when his actual demand is \hat{X}_i . We omit the membership fee r_i in the proof as it is irrelevant to firms' reporting decisions. Then, under the fixed-rate payment scheme, we have

$$h_i(\check{X}_i, \hat{X}_i) = p(\hat{X}_i - A_i(\check{X}_i))^+ + w_i^F A_i(\check{X}_i).$$

Under the allocation policy, $A_i(\check{X}_i) \leq \check{X}_i$ almost surely. Then,

$$h_i(\hat{X}_i, \hat{X}_i) = p\hat{X}_i - (p - w_i^F)A_i(\hat{X}_i),$$

which represents Firm i 's cost if he truthfully reports \hat{X}_i .

We show by contradiction that Firm i will be (weakly) worse off either overreporting or underreporting \hat{X}_i .

No Overreporting. Suppose that Firm i overreports, i.e., $\check{X}_i > \hat{X}_i$. Under our allocation policy, there are two possible scenarios: either $A_i(\check{X}_i) = \check{X}_i$ or $A_i(\check{X}_i) < \check{X}_i$.

- **Scenario (i):** $A_i(\check{X}_i) = \check{X}_i$. Under this scenario, the available workforce is sufficient to fulfill the request \check{X}_i . We have $h_i(\check{X}_i, \hat{X}_i) = w_i^F \check{X}_i$. If truthfully reporting \hat{X}_i instead, Firm i will be fully satisfied as well, i.e., $A_i(\hat{X}_i) = \hat{X}_i$, as $\hat{X}_i < \check{X}_i$. So, $h_i(\hat{X}_i, \hat{X}_i) = w_i^F \hat{X}_i \leq h_i(\check{X}_i, \hat{X}_i) = w_i^F \check{X}_i$ provided $w_i^F \geq 0$. That is, truthful reporting is optimal under Scenario (i).

- **Scenario (ii):** $A_i(\tilde{X}_i) < \hat{X}_i$. Under this scenario, Firm i 's request is not fully satisfied, implying that the available workforce is equal to $A_i(\tilde{X}_i)$. Depending on whether $A_i(\tilde{X}_i)$ is sufficient to cover the true demand \hat{X}_i , we have

$$h_i(\tilde{X}_i, \hat{X}_i) = \begin{cases} w_i^F A_i(\tilde{X}_i), & \text{if } \hat{X}_i \leq A_i(\tilde{X}_i) < \tilde{X}_i, \\ p\hat{X}_i - (p - w_i^F)A_i(\tilde{X}_i), & \text{if } \hat{X}_i > A_i(\tilde{X}_i). \end{cases}$$

In the first case, the available workforce $A_i(\tilde{X}_i)$ is sufficient to cover true demand \hat{X}_i , and so truthful reporting leads to $A_i(\hat{X}_i) = \hat{X}_i$. Thus, $h_i(\hat{X}_i, \hat{X}_i) = w_i^F \hat{X}_i \leq h_i(\tilde{X}_i, \hat{X}_i) = w_i^F A_i(\tilde{X}_i)$, i.e., the firm is better off by truthful reporting.

In the second case, the available workforce $A_i(\tilde{X}_i)$ is less than true demand \hat{X}_i . By truthful reporting, Firm i will receive the same amount as it would by overreporting, i.e., $A_i(\hat{X}_i) = A_i(\tilde{X}_i)$. Hence, $h_i(\hat{X}_i, \hat{X}_i) = h_i(\tilde{X}_i, \hat{X}_i)$ and the firm is indifferent between truthful reporting and overreporting.

Summarizing the above discussion, on all sample paths, Firm i 's cost under truthful reporting is lower than or equal to that under overreporting.

No Underreporting. Suppose Firm i underreports, i.e., $\tilde{X}_i < \hat{X}_i$. Then, we must have $A_i(\tilde{X}_i) \leq \tilde{X}_i < \hat{X}_i$ and so $h_i(\tilde{X}_i, \hat{X}_i) = p\hat{X}_i - (p - w_i^F)A_i(\tilde{X}_i)$. Truthful reporting will weakly increase its allocation, i.e., $A_i(\hat{X}_i) \geq A_i(\tilde{X}_i)$ such that $h_i(\hat{X}_i, \hat{X}_i) = p\hat{X}_i - (p - w_i^F)A_i(\hat{X}_i) \leq h_i(\tilde{X}_i, \hat{X}_i) = p\hat{X}_i - (p - w_i^F)A_i(\tilde{X}_i)$. Therefore, the firm has no incentive to underreport.

A.13 Proof of Proposition 3

Lemma 5 guarantees that each firm truthfully reports the actual vacancy $(X_i - Q_i)^+$. In what follows, we verify that, provided each firm truthfully reports his demand, the proposed mechanism induces the system-wide optimal solution as a Nash equilibrium.

Individual Firms. Given that the platform chooses the system-wide optimal η^* and allocates the workforce per the TFRB policy $\mathbf{A}^*(\cdot)$, Firm i 's expected cost function is adapted to

$$\begin{aligned} h_i(Q_i | \mathbf{Q}_{-i}, \mathbf{A}^*(\cdot), \eta^*) &= w_i^F \mathbb{E} [A_i^*((\mathbf{X} - \mathbf{Q})^+, S^*)] + p \mathbb{E} \left[((X_i - Q_i)^+)^+ - A_i^*((\mathbf{X} - \mathbf{Q})^+, S^*) \right] \\ &\quad + cQ_i + r_i \\ &= p(1 - m_i) \mathbb{E} [(X_i - Q_i)^+] + cQ_i + r_i, \end{aligned}$$

where the equality holds due to $\mathbb{E}[A_i^*] = \beta_i^* \mathbb{E}[(X_i - Q_i)^+]$ (by Lemma 1). As such, the expected cost function h_i remains the same as in Section 5.2.2, and so choosing m_i per (14) ensures that Firm i will set $Q_i = Q_i^*$.

Platform Operator. Under the proposed mechanism, the platform operator will in total receive an income $w_i^F A_i + w_i^C(A_i, \hat{X}_i) - r_0$ if allocating A_i to Firm i . That is, her expected profit remains the same as in Section 5.2.1 except paying an extra amount nr_0 to the third party. Since the extra payment is constant, the platform's incentive is still aligned with the system when deciding on η and the allocation rule. Thus, given all firms choosing $Q_i = Q_i^*$, it is optimal for the platform to choose the system-wide

optimal η^* and allocate according to the TFRB policy.

In summary, the system-wide optimal solution is a Nash equilibrium under the proposed mechanism. To ensure participation, we modify (16) to (20) to take into account the extra amount nr_0 the platform pays to the third party. As $\mathbb{E}[w_i^C(A_i, \hat{X}_i)] = 0$ under the TFRB policy, the third party's expected payoff is nr_0 .

B Discussion on Price-Only Contracts

We have shown that the staffing system can be coordinated with a three-parts payment scheme. One may wonder what happens under our context if we use a simple price-only contract under which each firm pays only a fixed rate, denoted by w^F , for each unit of allocated on-demand labor.

Under the price-only contract, the platform maximizes the following expected profit by choosing η and allocation policy $\mathbf{A}(\cdot) \in \mathcal{A}$:

$$\Pi(\eta, \mathbf{A}(\cdot)|\mathbf{Q}) = w^F \mathbb{E} \left[\min \left\{ \sum_{i \in I} (X_i - Q_i)^+, NG(\eta) \right\} \right] - \eta NG(\eta). \quad (\text{B.8})$$

Given other players' decisions, Firm i minimizes his cost function below:

$$\begin{aligned} h_i(Q_i|\mathbf{Q}_{-i}, \mathbf{A}(\cdot), \eta) &= p \mathbb{E} \left[((X_i - Q_i)^+ - A_i)^+ + w^F A_i \right] + cQ_i \\ &= p \mathbb{E} [(X_i - Q_i)^+] - (p - w^F) \mathbb{E}[A_i] + cQ_i, \end{aligned} \quad (\text{B.9})$$

where A_i depends on all firms' job vacancies $(\mathbf{X} - \mathbf{Q})^+$ and on-demand staffing level $S = NG(\eta)$. Unlike the case under our proposed contract, allocation A_i is not canceled out in Firm i 's problem. Thus, the platform operator can potentially manipulate the allocation rule to its own interest. As the allocation rule can be almost arbitrary, the problem of finding a Nash equilibrium for the game appears not very well-defined.

In the literature on inventory allocation, to bypass the above issue, existing papers assume that the platform can commit to a prespecified allocation rule and then analyze a noncooperative game under the given allocation rule (e.g., Cachon and Lariviere, 1999; Netessine and Rudi, 2006).

Hence, we follow a similar approach as adopted in the literature to analyze our staffing system under the price-only contract. Nonetheless, we emphasize that for all the results presented in our main paper to hold, we need *not* assume any commitment power for the platform.

Following the extant literature, we consider a widely studied allocation rule—the relaxed linear allocation (RLA) rule, which is defined as

$$A_i^{lin}((\mathbf{X} - \mathbf{Q})^+, S) = \min \left\{ (X_i - Q_i)^+, (X_i - Q_i)^+ - \frac{1}{n} \left(\sum_{k=1}^n (X_k - Q_k)^+ - S \right) \right\}, \forall i \in I. \quad (\text{B.10})$$

Under the RLA rule, each firm receives his request minus a common deduction. The literature has demonstrated the merits of the above allocation rule; we refer to interested readers to Cachon and

Lariviere (1999) and Netessine and Rudi (2006) for the detailed justifications of the RLA rule. In particular, the RLA ensures the existence of a pure-strategy Nash equilibrium in their problems, which is also true for our problem.²³

Under any given price-only contract with $w^F < p$ and the RLA rule, Firm i 's cost function is reduced to

$$\begin{aligned} \min_{Q_i \geq 0} h_i(Q_i | \mathbf{Q}_{-i}, \mathbf{A}^{lin}(\cdot), S) &= p \mathbb{E} [(X_i - Q_i)^+] - (p - w^F) \mathbb{E} [A_i^{lin}((\mathbf{X} - \mathbf{Q})^+, S)] + cQ_i \\ &= w^L \mathbb{E} [(X_i - Q_i)^+] + (p - w^L) \mathbb{E} \left[\left(\frac{1}{n} \sum_{k=1}^n (X_k - Q_k)^+ - S \right)^+ \right] + cQ_i. \end{aligned} \quad (\text{B.11})$$

On the other hand, the platform operator maximizes (B.8) by choosing η , while committing the RLA rule given by (B.10). It is convenient to change the platform operator's decision variable from η to on-demand staffing level S . The platform's problem is thus given by

$$\max_{S \geq 0} \Pi(S | \mathbf{Q}) = w^L \mathbb{E} \left[\min \left\{ \sum_{i \in I} (X_i - Q_i)^+, S \right\} \right] - SG^{-1} \left(\frac{S}{N} \right) \quad (\text{B.12})$$

Hence, the platform and n firms engage in a simultaneous-move game in which the platform operator chooses S and Firm i decides on Q_i .

Proposition B.1. Given any price-only contract with $\underline{\tau} < w^F < p$, under the RLA rule, a pure-strategy Nash equilibrium (S^N, \mathbf{Q}^N) exists. The game is submodular in the sense that the best response functions $Q_i^{br}(\mathbf{Q}_{-i}, S)$ for all $i \in I$ and $S^{br}(\mathbf{Q})$ are decreasing. Moreover, any interior Nash equilibrium (\mathbf{Q}^N, S^N) must satisfy the following set of optimality conditions:

$$\begin{aligned} w^L \bar{F}_i(Q_i^N) + \frac{p - w^L}{n} \mathbb{P} \left(S^N < \frac{1}{n} \sum_{k=1}^n (X_k - Q_k^N)^+, Q_i^N < X_i \right) &= c, \forall i \in I, \\ w^L \mathbb{P} \left(S^N < \sum_{k=1}^n (X_k - Q_k^N)^+ \right) &= G^{-1} \left(\frac{S^N}{N} \right) + \frac{S^N}{N} \frac{1}{g \left(G^{-1} \left(\frac{S^N}{N} \right) \right)}. \end{aligned}$$

Proof. The existence of a pure-strategy Nash equilibrium follows immediately by the convexity of h_i and concavity of Π . The submodularity of the game can be proved by checking the cross partial derivatives. Notice that $\partial^2 h_i / (\partial Q_i \partial S) \geq 0$ because $\partial h_i / \partial S = -(p - w^L) \mathbb{P}(S < (1/n) \sum_{k=1}^n (X_k - Q_k)^+)$, which is increasing in Q_i , and $\partial^2 h_i / (\partial Q_i \partial Q_j) \geq 0$ as $\partial h_i / \partial Q_j = -(p - w^L) / n \mathbb{P}(S < \frac{1}{n} \sum_{k=1}^n (X_k - Q_k)^+, Q_j < X_j)$, which is increasing in Q_i . By the convexity of h_i in Q_i , the minimizer $Q_i^{br}(\mathbf{Q}_{-i}, S)$ is decreasing in S and \mathbf{Q}_{-i} . Furthermore, it follows that $\partial^2 \Pi / (\partial S \partial Q_i) \leq 0$ by observing that $\partial \Pi / \partial S = w^L \mathbb{P}(S < \sum_{k=1}^n (X_k - Q_k)^+) - G^{-1}(S/N) + S/(Ng(G^{-1}(S/N)))$ is decreasing in each Q_i . As Π is concave in S , the maximizer $S^{br}(\mathbf{Q})$ is decreasing in each Q_i . Optimality conditions follow immediately from the first-order conditions. \square

The above proposition implies that the staffing decisions of the $(n + 1)$ players in the game are

²³Because of the term $(p - w^F) \mathbb{E}[A_i]$ in (B.9), each firm's cost function h_i may not be quasiconvex in general. Under the RLA rule, h_i is convex in Q_i and therefore, the existence of a Nash equilibrium is guaranteed.

strategic substitutes: As a firm or the platform increases its staffing level, the other players will choose a lower staffing level in equilibrium. This result contrasts with the coordinated system under our proposed mechanism in which the n firms' problems are decoupled into n independent problems.

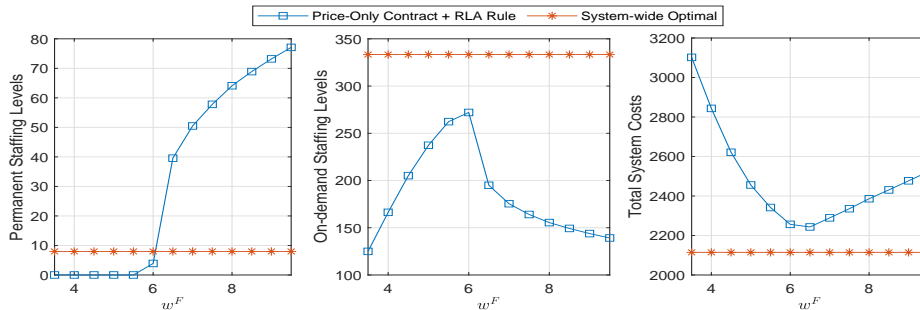


Figure 2: Price-only contract (with the RLA rule) versus the system-wide optimality

We consider a numerical example to illustrate how the equilibrium under the price-only contract and RLA allocation can deviate from the system-wide optimal solution.

For illustrative purposes, we consider four firms ($n = 4$) with identical and independently distributed demands. We can thus focus on the equilibrium symmetric to the four firms (i.e., the equilibrium where the Q_i^N 's are equal). Their demand distribution is assumed to be Gamma with mean = 100 and coefficient of variation = 0.5. Other parameters are set as $p = 10$, $c = 6$, $N = 1000$, and G is a uniform distribution between $[2, 8]$. With the fixed rate w^F varied from 3.5 to 9.5 with an increment of 0.5, we compute the equilibrium (S^N, \mathbf{Q}^N) under the price-only contract and RLA rule and compare it against the system-wide optimal solution. Results are plotted in Figure 2. Because of the gaming effect among the players, the staffing levels under the price-only contract and RLA rule can be substantially different from the system-wide optimal level. In particular, we observe that with price-only contracts, the platform will choose an on-demand staffing level lower than the system-wide optimal (see the middle panel of Figure 2), suggesting that the price-only contract does not provide sufficient incentive for the decentralized system to make the best use of on-demand staffing. As a consequence, the system cost under the price-only contract and RLA rule can be much higher than that under a coordinated system (see right panel of Figure 2).

C Effect of Fill Rate Refinement

C.1 Effect of Fill Rate Refinement (Algorithm 2)

Algorithm 2 presented in the main paper describes an approach to refine the optimal target fill rates such that the target fill rates offered to each firm are as close as possible. To demonstrate the power of Algorithm 2, for every instance in our numerical experiment reported in Section 6, we obtain two sets of fill rates by solving problem (8) with a greedy method based on its polymatroid structure (Edmonds, 1970) and by applying Algorithm 2, respectively. (For a detailed description of our numerical setup, see Section 6.2.) For each instance and each method, we calculate the maximum difference between any pair

of fill rates: $\max_{i,j} |\beta_i^* - \beta_j^*|$. As reported in Table 1, on average, Algorithm 2 can reduce the maximum difference in fill rates from 0.435 to 0.035 compared to the greedy method. Furthermore, in 83.1% of tested instances, the refined fill rates are nearly equal (with $\max_{i,j} |\beta_i^* - \beta_j^*|$ less than 10^{-6}).

Table 1: Effect of Fill Rate Refinement (Algorithm 2)

Marginal dis- tribution	Demand share of Firm 1	$\max \beta_i^* - \beta_j^* $ by greedy	$\max \beta_i^* - \beta_j^* $ by Algorithm 2	% of Instances with $\max \beta_i^* - \beta_j^* <$ 10^{-6}
Gamma	0.05	0.682	0.000	100.0%
	0.1	0.659	0.000	100.0%
	0.25	0.539	0.000	100.0%
	0.4	0.408	0.000	100.0%
	0.6	0.277	0.030	83.3%
	0.9	0.216	0.145	44.4%
	0.95	0.212	0.151	33.3%
LogNormal	0.05	0.672	0.000	100.0%
	0.1	0.641	0.000	100.0%
	0.25	0.515	0.000	100.0%
	0.4	0.387	0.000	100.0%
	0.6	0.263	0.032	77.8%
	0.9	0.205	0.131	50.0%
	0.95	0.200	0.136	44.4%
Uniform	0.05	0.588	0.000	100.0%
	0.1	0.590	0.000	100.0%
	0.25	0.598	0.000	100.0%
	0.4	0.586	0.000	100.0%
	0.6	0.444	0.000	100.0%
	0.9	0.197	0.041	66.7%
	0.95	0.181	0.046	66.7%
Mix	0.05	0.650	0.010	91.7%
	0.1	0.589	0.002	96.7%
	0.25	0.589	0.000	100.0%
	0.4	0.465	0.003	96.7%
	0.6	0.379	0.054	70.0%
	0.9	0.230	0.078	56.7%
	0.95	0.214	0.111	48.3%
Average		0.435	0.035	83.1%

C.2 Suboptimality of Equal Target Fill Rates

Example 1. (SUBOPTIMALITY OF EQUAL TARGET FILL RATES) Consider two firms. Firm 1's demand follows a two-point distribution as follows:

$$X_1 = \begin{cases} 5, & \text{with probability } 1/5, \\ 0, & \text{with probability } 4/5. \end{cases} \quad (\text{C.13})$$

Firm 2's demand follows a uniform distribution between 0 and 2, i.e., $X_2 \sim \text{Uniform}(0, 2)$. Assume that X_1 and X_2 are independent. Suppose the model parameters are such that $S^* = 1$ and $Q_1^* = Q_2^* = 0$. We have $\mathbb{E}[(X_1 - Q_1^*)^+] = \mathbb{E}[(X_2 - Q_2^*)^+] = 1$, $\mathbb{E}[\min\{S^*, X_1\}] = 1/5$, $\mathbb{E}[\min\{S^*, X_2\}] = 3/4$, and

$\mathbb{E}[\min\{S^*, X_1 + X_2\}] = 4/5$. Problem (8) reduces to

$$\begin{aligned} \max \quad & \beta_1 + \beta_2 \\ \text{s.t.} \quad & \beta_1 + \beta_2 \leq 4/5 \\ & \beta_1 \leq 1/5 \\ & \beta_2 \leq 3/4 \end{aligned}$$

The optimal objective value is $4/5$, which can be attained by an optimal solution $(\beta_1^*, \beta_2^*) = (1/5, 3/5)$. This solution is also one in which the minimum fill rate is maximized, as β_1^* has achieved its upper bound $1/5$. However, if we restrict to a firm-independent fill rate (i.e., $\beta_1 = \beta_2$), the solution will be $(1/5, 1/5)$, which leads to an objective value of $2/5$, much lower than the actual optimal value $4/5$.

D A Sufficient Condition for $w_i^F \geq 0$

Proposition D.2. If X_i 's are exchangeable and exclusive use of on-demand staffing is system-wide optimal, provided that the aggregate demand $\sum_{i \in I} X_i$ is a log-concave random variable with no probability mass at zero, then $\beta_i^* \geq m_i$ and $w_i^F \geq 0$ for all $i \in I$.

Proof. For exchangeable demands, Proposition 1 applies, and therefore we have

$$\beta_i^* = \frac{\mathbb{E}[\min\{S^*, \sum_{i \in I} X_i\}]}{n \mathbb{E}[X_i]} = 1 - \frac{\mathbb{E}[(\sum_{i \in I} X_i - S^*)^+]}{n \mathbb{E}[X_i]}.$$

On the other hand, $m_i = 1 - c/p$ as $Q_i^* = 0$ for all $i \in I$. Thus, it suffices to show

$$\frac{c}{p} \geq \frac{\mathbb{E}[(\sum_{i \in I} X_i - S^*)^+]}{n \mathbb{E}[X_i]}.$$

Claim 4. If Y is a nonnegative random variable with a continuous log-concave distribution F and $F(0) = 0$, then for any given scalar $S \geq 0$, we have $\mathbb{P}(Y > S) \geq \mathbb{E}[(Y - S)^+]/\mathbb{E}[Y]$.

Proof of Claim 4. To prove this claim, it is convenient to avail of the following result, which is proved in Proposition 1 of Heckman and Honore (1990): *If Y is a log-concave random variable, then $0 \leq \partial \mathbb{E}[Y|Y > s]/\partial s \leq 1$.*

It follows that $-1 \leq \partial \mathbb{E}[Y - s|Y > s]/\partial s \leq 0$, so $\mathbb{E}[Y - s|Y > s]$ is decreasing in s . Thus, since $F(0) = 0$, we have

$$\mathbb{E}[Y] = \mathbb{E}[Y|Y > 0] \geq \mathbb{E}[Y - S|Y > S], \text{ for any scalar } S.$$

As $\mathbb{E}[(Y - S)^+] = \mathbb{E}[Y - S|Y > S]\mathbb{P}(Y > S)$, multiplying $\mathbb{P}(Y > S)$ on both sides of $\mathbb{E}[Y] \geq \mathbb{E}[Y - S|Y > S]$ and rearranging the terms lead to the desired inequality. \square

Following the assumption on $\sum_{i \in I} X_i$ and the above claim, we have

$$\mathbb{P}\left(\sum_{i \in I} X_i > S^*\right) \geq \frac{\mathbb{E}[(\sum_{i \in I} X_i - S^*)^+]}{n \mathbb{E}[X_i]}.$$

By the optimality condition given in Theorem 1, when on-demand staffing is exclusively used in the system-wide optimal solution, we have $c \geq \phi^o = p\mathbb{P}(\sum_{i \in I} X_i > S^*)$, i.e., $c/p \geq \mathbb{P}(\sum_{i \in I} X_i > S^*)$. These two inequalities together imply that $c/p \geq \mathbb{E}[(\sum_{i \in I} X_i - S^*)^+]/(n \mathbb{E}[X_i])$, which completes the proof. \square

E When On-demand Workers Are Not Infinitesimal

We have assumed that on-demand workers are infinitesimal such that the amount of labor supply from the available on-demand workforce induced by a given compensation rate η is a deterministic number, $NG(\eta)$, where N is the total amount of potential labor supply. In fact, our main results on system coordination do *not* rely on this assumption. We can relax the assumption of infinitesimal workers to a finite number of potential on-demand workers, each of whom offers K units of labor supply and consequently S/K follows a binomial distribution with parameters N/K and $G(\eta)$, where K is a constant such that N/K is an integer. As K goes to zero, the system becomes close to the one with infinitesimal workers.

The analysis of system coordination hinges on showing that the objective function in the problem faced by the self-interested platform operator coincides with the system-wide objective function, $C(\mathbf{Q}, \eta)$ in (5). The analysis remains unchanged by treating S in $\mathbb{E}[(\sum_{i \in I} (X_i - Q_i)^+ - S)^+]$ as a random variable and taking the expectation over both X_i 's and S . Then the two objective functions are still aligned with each other, and consequently the platform operator will be incentivized to set the system-wide optimal compensation η^* . For the allocation policy and firm's incentive compatibility analysis, we can also treat S as a random variable in the constraints of problem (8) and in optimal target fill rate computation (see Appendix A.4.3), and take the expectations over both X_i 's and S . Therefore, by setting contract parameter m_i 's according to Lemma 3, the individual firm's incentive will still be aligned with the centralized system and be willing to choose the system-wide optimal permanent staffing level \mathbf{Q}^* . To summarize, the proposed contract scheme and allocation policy can still lead to system coordination in the decentralized system under stochastic labor supply S , i.e., Theorem 2 still holds. Therefore, all insights on the contract design remain valid.

The only issue arises when one is interested in numerically calculating some contract parameters, in particular, m_i 's. Note that for another set of contract parameters, β_i 's, one can still easily compute using Lemma 1. The calculation of m_i 's with (14) in Lemma 3, however, requires the value of \mathbf{Q}^* , the system-wide optimal permanent staffing level, which is the solution to the system-wide problem with objective function (5). When S becomes a random variable such that S/K follows a binomial distribution with parameter N/K and $G(\eta)$, the system-wide problem is generally not jointly convex in η and \mathbf{Q} . One solution approach is to approximate the objective function using the infinitesimal-worker assumption by replacing S with its mean $NG(\eta)$. Note that by Jensen's inequality, the approximated objective function is a lower bound to the original one. For any given η , we can derive a bound on the gap between the

approximation and original object function as follows:

$$\begin{aligned}
& \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i)^+ - S \right)^+ \right] - \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i)^+ - NG(\eta) \right)^+ \right] \\
& \leq \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i)^+ - NG(\eta) \right)^+ + (NG(\eta) - S)^+ \right] - \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i)^+ - NG(\eta) \right)^+ \right] \\
& = \mathbb{E} \left[(NG(\eta) - S)^+ \right] = K \mathbb{E} \left[\left(\frac{N}{K} G(\eta) - \frac{S}{K} \right)^+ \right] \\
& \leq K \cdot \frac{1}{2} \left\{ \frac{N}{K} G(\eta) - \mathbb{E} \left[\frac{S}{K} \right] + \sqrt{\left(\frac{N}{K} G(\eta) - \mathbb{E} \left[\frac{S}{K} \right] \right)^2 + \text{Var} \left(\frac{S}{K} \right)} \right\} \\
& = \frac{K}{2} \sqrt{\frac{N}{K} G(\eta) [1 - G(\eta)]} = \frac{1}{2} \sqrt{K N G(\eta) [1 - G(\eta)]},
\end{aligned}$$

where the second inequality in the third last row follows the well-known distribution-free bound by Scarf (1958), and $\text{Var}(S/K)$ denotes the variance of random variable S/K . We can further develop a bound independent of η by replacing $G(\eta)$ with $1/2$ and obtain $\sqrt{KN}/4$. The bound will converge to zero at a rate in the order of square root of K as K approaches zero. Under such approximation, the system-wide optimal solution remains the same as stated in Theorem 1, as the approximated objective function is the same as the one in the base model under the assumption of infinitesimal workers.

F Workers Being Paid Only When They Get Actual Jobs

Thus far, we have assumed that workers are paid as long as they are on call (i.e., willing to work). An alternative wage scheme—in which workers are paid only when they complete an actual job—has been nicely modeled by Cachon et al. (2017) for a platform facing uncertain demand. In this appendix, we will show that our coordination mechanism is still valid when this alternative wage scheme is adopted in a decentralized system consisting of a platform operator and multiple self-interested employers.

Similar to Cachon et al. (2017), let θ denote the fraction of labor supply offered by participating workers used to fill job positions. $\theta = 1$ if total demand exceeds the total labor supply, whereas $\theta < 1$ if supply is rationed and there are participating workers who are not assigned to any job positions. The N temporary workers on the platform would therefore decide to participate according to rational expectations. The value of θ in equilibrium is determined by the recursive relation:

$$\theta = \begin{cases} 1, & \text{if } NG(\eta) \leq \sum_{i \in I} (X_i - Q_i)^+, \\ \frac{\sum_{i \in I} (X_i - Q_i)^+}{NG(\theta\eta)}, & \text{if } NG(\eta) > \sum_{i \in I} (X_i - Q_i)^+. \end{cases} \quad (\text{F.14})$$

The system-wide optimal solution can thus be obtained by solving the following problem.

$$\begin{aligned}
\min_{\mathbf{Q}, \eta} \quad & p \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i)^+ - NG(\theta\eta) \right)^+ \right] + \eta NG(\theta\eta) + c \sum_{i \in I} Q_i \\
\text{s.t.} \quad & \theta \text{ and } \eta \text{ satisfy (F.14)} \\
& Q_i \geq 0, \forall i \in I
\end{aligned} \tag{F.15}$$

Computing exact optimal compensation rates becomes more challenging, as (F.14) should be satisfied on any sample path. Nevertheless, it can be shown without explicitly solving problem (F.15) that the mechanism proposed in Section 5 can still induce system optimality in decentralized systems. When the staffing level of Firm i is fixed as Q_i^* , following the same derivation in Section 5.2, the platform operator's problem can be reduced to

$$\begin{aligned}
\max_{\eta} \quad & -\eta NG(\theta\eta) - p \mathbb{E} \left[\left(\sum_{i \in I} (X_i - Q_i^*)^+ - NG(\theta\eta) \right)^+ \right] \\
& - p \sum_{i \in I} (1 - m_i^*) \mathbb{E} \left[(X_i - Q_i^*)^+ \right] + \sum_{i \in I} r_i \\
\text{s.t.} \quad & \theta \text{ and } \eta \text{ satisfy (F.14)}
\end{aligned}$$

which is aligned with the objective of the entire system. Thus, the platform will set η as the system-wide optimal level and assign the on-demand workforce using the desired allocation rule. On the other hand, individual employers' problems remain intact under this alternative wage scheme, since each has been guaranteed a given fill rate no matter how temporary workers are being paid. Therefore, the proposed coordination mechanism is still applicable. To implement it, one can first numerically find the system-wide optimal η^* and Q_i^* by solving problem (F.15), and then construct the incentive contracts as in Section 5.2.

G Alternative Allocation under Verifiable Job Vacancies

The TFRB policy described in Algorithm 1 allocates full requested amounts to the firms at the top of the priority list until the total amount of on-demand workforce is depleted or every firm's requested demand is satisfied. As shown in Section 7, such allocation induces true-telling behavior from the firms when the demand is not verifiable. One key behind the results is the full demand allocation, which will hurt any firm that overreports his demand. In practice, when the supply of on-demand workforce is limited, such a full demand allocation rule will lead to extreme situations in which some firms' requested demands are fully satisfied while others receive zero on-demand workers. The platform may want to make the allocation "fairer"²⁴ by assigning on-demand workers to as many firms as possible. We present an

²⁴Note that fairness is a complicated issue with various versions of definitions. We only use the term very loosely here to describe the allocation that allows more firms to enjoy on-demand workers. It is beyond the scope of the current paper to discuss fairness in the on-demand workforce allocation, which could be a potential future research direction.

alternative allocation policy by adjusting the allocation mechanism in the TFRB policy and show that such an adjusted allocation mechanism could still coordinate the system when the demands are verifiable.

The setup of the alternative allocation is the same as the TFRB policy presented in Algorithm 1, and the adjustments are made to Step 1b and Step 3 during the demand allocation. The priority lists are derived in the same manner as stated in Step 1b and Step 2. In Step 1b, instead of allocating full demand, the alternative approach allocates $\beta_{[i]}\hat{X}_{[i]}(t)$ to the i^{th} firm on the priority list and then moves to the next firm on the list. This ensures the target fill rates of the firms are met in the single period as much as possible. The two scenarios in Step 1b are revised to:

- Scenario 1: $\sum_{i \in I} \beta_{[i]}\hat{X}_{[i]}(t) > S$, i.e., the workforce is not sufficient to satisfy the target fill rates of all firms. Define $i' = \min\{j : \sum_{i=1}^j \beta_{[i]}\hat{X}_{[i]}(t) > S\}$. Then the allocation stops at Firm $[i']$ with $A_{[i]}(t) = \beta_{[i]}\hat{X}_{[i]}(t)$, for all $i < i'$, and Firm $[i']$ receives the remaining amount of workforce, i.e., $A_{[i']}(t) = S - \sum_{j=1}^{i'-1} \hat{X}_{[j]}(t)$. The rest firms receive zero allocation.
- Scenario 2: $\sum_{i \in I} \hat{X}_{[i]}(t) \leq S$, i.e., the workforce is sufficient to satisfy the target fill rates of all firms. Any remaining amount of workforce can be allocated arbitrarily among the firms. Then $A_{[i]}(t) \geq \beta_{[i]}\hat{X}_{[i]}(t)$, for all $i \in I$.

Similarly, Step 3 in Algorithm 1 is revised accordingly to: “Allocate $\beta_i\hat{X}_i$ amount of workforce to Firm i following the priority list obtained in Step 2 until reaching the end of the list, or the available workforce is depleted; remaining on-demand workers, if any, can be arbitrarily allocated to satisfy unmet demands.” Note that in such an allocation scenario, the remaining on-demand workers can still be utilized as the priority allocation stage only fulfills $\beta_i\hat{X}_i$ amount of demands from Firm i and Firm i 's remaining demand $(1 - \beta_i)\hat{X}_i$ is still unsatisfied. This is different from our base model, in which the priority allocation stage fulfills the full demand requests from the firms. We refer to the revised allocation policy as Alternative Target Fill Rate-Based (ATFRB) policy.

To see that Lemma 1 holds with the ATFRB policy, it suffices to verify Step 1 in the proof of Lemma 1. To this end, the only change happens in the proof of Claim 2.

Proof of Claim 2 Under the ATFRB Policy. There are two possible cases when the allocation process stops, which we discuss separately. In the first case, the allocation stops at the $(k + 1)$ th firm in the priority list. According to the ATFRB policy, this will happen only when the remaining capacity, $S - \sum_{i=1}^k A_i(t)$, is less than 0, i.e., $A_{k+1}(t) = 0$ if $\sum_{i=1}^{k-1} A_i(t) \leq S$ and $\sum_{i=1}^k A_i(t) \geq S$. For any $i \leq \min\{k, m\}$, $A_i(t) = \beta_i\hat{X}_i(t)$. Then we have $\mathbb{E}[\sum_{i=1}^j A_i(t)] = \sum_{i=1}^j \beta_i \mathbb{E}[\hat{X}_i]$, for all $j = 1, 2, \dots, \min\{k, m\}$. If $m \leq k$, the claim is proved. If $m > k$, $\sum_{i=1}^j A_i(t) \geq S \geq \min\{S, \sum_{i=1}^j \hat{X}_i(t)\}$, for all $j = k + 1, k + 2, \dots, m$. Taking the expectations on both sides of the above inequality and applying the constraints in problem (8), we have

$$\mathbb{E} \left[\sum_{i=1}^j A_i(t) \right] \geq \mathbb{E} \left[\min \left\{ S, \sum_{i=1}^j \hat{X}_i(t) \right\} \right] \geq \sum_{i=1}^j \beta_i \mathbb{E}[\hat{X}_i], \forall j = k + 1, k + 2, \dots, m,$$

which implies the inequalities stated in the claim.

In the second case, the allocation process does not stop after going through all of the firms on the priority list once. According to the ATFRB policy, it is obvious that for any Firm i , $A_i(t) \geq \beta_i \hat{X}_i(t)$. The claim follows immediately after taking the expectations. \square

Finally, the system coordination result in Theorem 2 follows from Lemma 1 that the ATFRB policy can also deliver the target fill rates to all firms in expectation, i.e., $\mathbb{E}[A_i^*] \geq \beta_i^* \mathbb{E}[(X_i - Q_i)^+]$, for all $i \in I$, and the ATFRB policy also satisfies the no-waste condition. These properties of the allocation policy are the key behind the individual incentive result established in Lemma 3. To summarize, as the ATFRB policy has the same properties as the TFRB policy when demands are verifiable, it can achieve system coordination under the same contract mechanism developed in Section 5.2.

References

- Blackwell, D. 1956. An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics* **6**(1) 1–8.
- Cachon, G. P., K. M. Daniels, R. Lobel. 2017. The role of surge pricing on a service platform with self-scheduling capacity. *Manufacturing & Service Operations Management* **19**(3) 368–384.
- Cachon, G. P., M. A. Lariviere. 1999. Capacity choice and allocation: Strategic behavior and supply chain performance. *Management Science* **45**(8) 1091–1108.
- Edmonds, J. 1970. Submodular functions, matroids, and certain polyhedra. *Combinatorial structures and their applications* 69–87.
- Heckman, James J, Bo E Honore. 1990. The empirical content of the roy model. *Econometrica: Journal of the Econometric Society* 1121–1149.
- Lu, Yingdong, David Yao. 2008. Linear optimization over a polymatroid with side constraints—scheduling queues and minimizing submodular functions. *arXiv:0804.1603 [math.OC]* .
- Netessine, S., N. Rudi. 2006. Supply chain choice on the internet. *Management Science* **52**(6) 844–864.
- Scarf, H. 1958. A min-max solution of an inventory problem. K. Arrow, S. Karlin, Scarf, H., eds., *Studies in The Mathematical Theory of Inventory and Production*. Stanford University Press, Redwood City, CA, 201–209.
- Shaked, M., J. G. Shanthikumar. 1997. Supermodular stochastic orders and positive dependence of random vectors. *Journal of Multivariate Analysis* **61**(1) 86–101.
- Topkis, Donald M. 1998. *Supermodularity and Complementarity*. Princeton University Press.
- Zhong, Y., Z. Zheng, M. C. Chou, C.-P. Teo. 2017. Resource pooling and allocation policies to deliver differentiated service. *Management Science* **64**(4) 1555–1573.