



# Automated Market Making for Goods with Perishable Utility

Chengqi Zang<sup>1,2</sup>, Gabriel P. Andrade<sup>1</sup> and Oğuzhan Ersoy<sup>1</sup>

<sup>1</sup>Gensyn AI, <sup>2</sup>The University of Tokyo

We study decentralized markets for goods whose utility perishes in time, with compute as a primary motivation. Recent advances in reproducible and verifiable execution allow jobs to pause, verify, and resume across heterogeneous hardware, which allow us to treat compute as time-indexed capacity rather than bespoke bundles. We design an automated market maker (AMM) that posts an hourly price as a concave function of load—the ratio of current demand to a “floor supply” (providers willing to work at a preset floor). This mechanism decouples price discovery from allocation and yields transparent, low-latency trading. We establish existence and uniqueness of equilibrium quotes and give conditions under which the equilibrium is admissible (i.e. active supply weakly exceeds demand). To align incentives, we pair a premium-sharing pool (base cost plus a pro-rata share of contemporaneous surplus) with a Cheapest-Feasible Matching (CFM) rule; under mild assumptions, providers optimally stake early and fully while truthfully report costs. Despite being simple and computationally efficient, we show that CFM attains bounded worst-case regret relative to an optimal benchmark.

## 1. Introduction

Two-sided markets offer a natural solution for many domains where interactions between distinct groups are beneficial but, due to these groups having their own externalities and constraints, a mediator is required to align incentives, optimize matching, and balance supply and demand (Rochet and Tirole, 2006; Jullien et al., 2021). Ride-sharing (Lai, 2021), fresh food delivery (Wang, 2022), and car rental (Hu and Zhou, 2015) platforms are well-known examples of two-sided markets at work, where these platforms serve the important role of managing prices and payouts while ensuring conditions are met (e.g. timeliness) for both sides to extract utility.

Having become one of the world’s most valuable commodities with rapidly growing demand (Data Bridge Market Research, 2024; Grand View Research, 2024), compute stands out as a compelling candidate for two-sided market design; however, the market for compute is characterized by sharp asymmetries and inefficiencies on both sides. Demand is often defined by heterogeneous and user-specific constraints—such as resource requirements, budget, deadlines, and seasonality—that make flexible matching essential. Supply, by contrast, often comes from large centralized providers that force users into rigid access models (e.g. on-demand, spot, or reserved instances) with complex tradeoffs between price, availability, and reliability—often resulting in substantial and rising long-term costs (Perez-Salazar et al., 2018; Furman and Diamant, 2020). At the same time, since privately held compute resources frequently suffer from chronic under-utilization and idle capacity while generally having prohibitive upfront costs (Gong et al., 2023), overall utilization of compute infrastructure remains surprisingly low.

These mismatches between dynamic, constraint-laden demand and inflexible, unevenly distributed supply reflect the absence of a truly open and liquid market where disparate providers or underutilized resources can easily participate. Instead, fragmented compute markets remain concentrated in a few large entities, who set opaque prices that reflect strategic or contractual considerations rather than transparent supply-and-demand dynamics. What is missing, therefore, is an efficient decentralized



marketplace for compute (Frick and Salmon, 2025), but prior work proposing such marketplaces has generally run into three persistent issues: Firstly, commonly utilized pricing mechanisms exacerbate the latency inherent to decentralization (e.g. auction mechanisms), exclude smaller participants, and fail to reflect real-time supply-demand dynamics (Bodoh-Creed et al., 2021; Wu et al., 2024). Secondly, standard approaches to optimal job matching introduce a fundamental computational complexity bottleneck (Leyton-Brown et al., 2002; Dobzinski and Vondrák, 2012) that many practical systems resolve with heuristics or batching, but this approach typically cannot guarantee incentive compatibility and is poorly suited for real-time trade execution with decentralized participants (Ashlagi et al., 2020; Lee, 2016). Lastly, decentralized or multi-provider settings introduce challenges of trust and incomplete information (e.g. providers misrepresenting their resource capacity) that disallow standard economic assumptions about participant reliability.

In this paper, we derive a decentralized two-sided market design that addresses these fundamental challenges. Unlike prior work, our framework assumes a setting where reproducible checkpointing and verification are possible due to recent advancements by Arun et al. (2025); reliable reproducibility allows us to pause, verify, and resume tasks across heterogeneous hardware without loss of correctness, which implies that matching no longer requires complicated constraints that ensure supply-side providers can accept demand-side tasks in their entirety (§2). Leveraging these assumptions allows us to treat compute as a time-bound, perishable commodity that facilitates a two-sided market design with dynamic and transparent price discovery. We formalize this idea by introducing a market structure wherein compute providers stake collateral to mint time-bounded assets and register availability windows for their resources, enabling efficient allocation while aligning incentives with demand-side users (§3). To incentivize truthful participation on the supply-side we introduce a pooling mechanism that treats these providers as liquidity suppliers (§4), rewarding them with both a reported base price and a share of market surplus while prioritizing competitively priced suppliers during matching—thereby improving utilization and discouraging rent-seeking behaviours.

Building on these ideas, we decouple pricing from matching which allows us to sidestep the intractability of combinatorial auctions while preserving efficiency and robustness in allocations; equilibrium prices are continuously computed by an algorithmic market maker, and matching is performed via a straightforward greedy algorithm. We establish the existence, uniqueness, and admissibility of equilibrium prices (§5) and we show that greedy matching achieves bounded regret relative to an online optimum in addition to being incentive-compatible (§6). Finally, we extend the framework to an incomplete-information setting by introducing slashing-based verification (Arun et al., 2025) and a racing mechanism (§7), which incentivizes truthful reporting and ensures timely task execution under minimal trust assumptions. Taken together, these components lay the foundation for a theoretically-founded model of decentralized compute markets grounded in practical trust and verification assumptions.

## 2. Background & Related Work

### 2.1. Two-Sided Platforms & Dynamic Pricing for Time-Coupled Goods

Two-sided market theory studies how intermediaries internalize cross-side externalities and set prices and matching rules across both sides (Rochet and Tirole, 2006; Jullien et al., 2021). Beyond static fees, platforms increasingly rely on dynamic pricing to guide participation, manage congestion, and steer growth over time (Rys and Sobolewski, 2020; Varma et al., 2023; Cao et al., 2025). In our setting, the traded object is time, thus prices are naturally time-indexed and allocation is a scheduling problem. Consequently, the platform’s role is to set intertemporal incentives that balance short-run liquidity with longer-run participation and reliability.

Dynamic pricing for perishable and time-sensitive assets provides core tools—models with finite horizons, stochastic arrivals, and myopic buyers underpin revenue and welfare analyses for perishables (Gallego



and Hu, 2014; Gisches et al., 2021; Feng and Gallego, 2000; Anjos et al., 2005). However, compared to physical perishables, compute time in decentralized markets is “sharper” since providers are individuals with idle devices over limited periods of time; their unused capacity vanishes immediately and value decays continuously as usability windows shrink. The closest macro-scale analog is electricity, where unconsumed generation is effectively wasted and time-varying prices can synchronize consumption with marginal cost to reduce peaks and improve welfare (Harding et al., 2023; Joskow and Wolfram, 2012; Sweeney et al., 2024). That said, modern power systems with renewables, storage, and curtailment complicate the “waste” analogy. Our compute setting must additionally handle heterogeneous devices, diverse job durations, and strategic behavior—features that call for pricing and allocation rules tailored to decentralized, heterogeneous supply.

## 2.2. Compute Markets & Auctions

Cloud markets combine spot (bid-based) and on-demand (availability premia) modalities, with significant work on user risk, substitutions, and pricing trade-offs (Kash et al., 2019; Hoy et al., 2016; Dierks and Seuken, 2022). Expressive mechanisms—such as combinatorial auctions and continuous double auctions—can capture complementarities and multi-attribute resources, but they face scalability and winner-determination hardness, particularly in high-frequency or real-time settings (Shang et al., 2010; Prasad et al., 2016; Peng et al., 2023; Leyton-Brown et al., 2002; Dobzinski and Vondrák, 2012). These computational and latency frictions help explain why many practical systems resort to batching or heuristics, often at the expense of incentive guarantees. Against this backdrop, designs that reduce dimensionality—so that clearing does not require large, latency-prone auctions—are attractive.

## 2.3. Reproducibility, Verification, & Simplified Matching

Recent advances in reproducible operators (RepOps), deterministic replay, and verifiable checkpointing make compute effectively fungible in time (Arun et al., 2025). RepOps standardizes numerical semantics and operator order, yielding bitwise-identical results across heterogeneous accelerators; canonical checkpoints with cryptographic commitments allow progress to be paused, verified, and resumed; and refereed delegation identifies the first divergent step at low cost. With these tools, feasible matching between supply-side and demand-side no longer hinges on bespoke multi-attribute bundles or bilateral trust since a scheduler simply needs to partition computational graphs into units that providers’ machines can handle: supply can be represented as intervals of verifiable compute time, and demand as deadline- and budget-constrained run lengths.

Before this reproducibility and cheap verification, cross-device migration risked numerical drift, pre-emption often wasted work, and correctness could not be established cheaply; hence prior markets leaned on expressive bids and combinatorial auctions despite NP-hard winner determination and latency concerns in online settings (Leyton-Brown et al., 2002; Dobzinski and Vondrák, 2012). Using this stack, allocation reduces to time-indexed matching rather than bundle selection. In this reduced setting, the assignment problem is naturally modeled as online bipartite matching: jobs arrive over time and can be matched to any provider whose available window covers their required hours. Classic RANKING/greedy approaches give robust guarantees and real-time scalability, with extensions for weights and priorities (Karp et al., 1990; Brodal et al., 2007; Huang et al., 2024; Gupta, 2024). These online-matching tools complement scheduling results from distributed ML and systems work that likewise frame cluster assignment as time-constrained matching (Bao et al., 2018; Zhang et al., 2025; Deng et al., 2021). This conceptual simplification—enabled by reproducibility and verification—motivates designs that decouple price discovery from allocation and execute trades in real time without auction bottlenecks.



### 3. Market Setting, Structure, and Pricing

Suppose we have partitioned sets of hardware into appropriate *tiers*, where a “tier” loosely defines a set of machines having comparable performance<sup>1</sup>. Each tier of compute can be thought of as its own market and will have its own price per hour set by the mechanism. Therefore, for simplicity, throughout the remainder of this paper, we will speak in terms of a specific tier of compute but the same will apply regardless of compute tier.

Throughout this section we rigorously define the market participants, pricing mechanism, and equilibrium concept used for the two-sided markets studied in this paper. At a high level, the demand-side are users with compute jobs (defined by inputs, computational graph, etc.) and the supply-side are compute providers with machines they are willing to temporarily make available for others’ jobs. Users come to the market with budgets, deadlines, and an estimate of how long their job will take. Providers place a stake depending on how long they are willing to “lock in” their machine and specify their operational cost, i.e. minimum price they are willing to work at. Prices are posted according to an algorithmic pricing function, where the current market price maps a notion of load to the per-period quote. Users see the current price and decide whether to start their job, whereas providers simply ensure their “locked in” machines are available whenever prices drop below their operational cost. Allocation is handled separately by a feasibility-aware greedy matcher introduced in §6.

Below we proceed in three steps: (i) Define participants and the period timeline (§3.1); (ii) Specify the automated pricing mechanism that posts prices relative to current market load (§3.2); and (iii) Collect the equilibrium-quote conditions (existence, uniqueness, admissibility) that the rest of the paper builds on (§3.3).

#### 3.1. Participants & Market Dynamics

We characterize participants in the two-sided market—providers and users—and how they interact at each period. Time in the market is modeled as being discrete, i.e.  $t \in \mathbb{N}$ . At the beginning of period  $t$  the market observes the set of currently staked providers  $\mathcal{S}^t$  and the set of active/pending jobs  $\mathcal{D}^t$ . Price  $P^t$  posted at  $t$  applies for the whole period  $t$  and is re-computed at  $t + 1$ .

Importantly, throughout the paper we make a *no outside option* assumption for market participants, which is a natural assumption in mechanism and market design literature (Krishna and Perry, 1998; Mookherjee, 2006; Figueroa and Skreta, 2007).

**Assumption 1.** Outside the network, a provider’s idle compute—available only during its idle window—has no salvage value if not matched within that window, yielding zero payoff; a user derives value from a job only if it is (partially) completed by her stated deadline, and otherwise is indifferent to not running it. Accordingly, both sides’ reservation utilities are normalized to zero.

This assumption isolates the core allocation problem by removing outside arbitrage opportunities that would otherwise distort incentives.

**Providers (Supply-side).** Provider  $s \in \mathcal{S}^t$  reports a per-hour cost  $\hat{c}_s \geq 0$  and a remaining *availability window*  $\tau_s(t) \in \mathbb{N}$ . Availability decays monotonically while staked:

$$\tau_s(t+1) = \max\{\tau_s(t) - 1, 0\},$$

---

<sup>1</sup>The precise definition of a “tier” is not especially important. In practice, it only matters for establishing expectations around the “least performant” set of machines in a market, which can influence certain parameter choices in real-world settings.



unless the provider restakes. The (unobserved) *true cost* is defined as  $c_s$ , which is the minimum price at which  $s$  would provide their compute.

A staked provider with  $\hat{c}_s \leq P^t$  is called *active* (i.e. eligible for matching at  $t$ ); otherwise they are called *dormant*. The set of active providers is further compartmentalized into *idle* and *assigned* suppliers depending on whether they are currently assigned to a job or not. When assigned a provider is always paid at least their  $\hat{c}_s$  but could get paid more as prices increase due to a *pool sharing mechanism* introduced in §4.1.

**Users (Demand-side).** A user's job  $d \in \mathcal{D}^t$  is a tuple  $(B_d, T_d, w_d)$  with budget  $B_d > 0$ , deadline  $T_d \in \mathbb{N}$ , and minimum viable run length  $\underline{w}_d \in \mathbb{N}$ . Given a posted price  $P^t$ , the user can purchase any integer number of hours

$$w \in \{0\} \cup \{\underline{w}_d, \dots, U_d(P^t)\}, \quad U_d(P) := \min\left\{\left\lfloor \frac{B_d}{P^t} \right\rfloor, T_d\right\}.$$

The job's value  $v_d(w)$  is assumed to be increasing and discretely concave (explored in §5.2), and users choose hours  $w_d(P^t)$  by solving

$$\max_{w \in \{0\} \cup \{\underline{w}_d, \dots, U_d(P^t)\}} v_d(w) - P^t w.$$

**Market Dynamics.** At each period  $t$ :

1. **Floor supply and total demand.** The *floor price*  $P_f$  is a pre-specified price<sup>2</sup> representing a reasonable lower bound for the “tier” of machine we expect in the market. Define the *floor supply* as the number of providers whose reported cost is below the floor price, i.e.

$$S_f^t := |\{s \in \mathcal{S}^t : \hat{c}_s \leq P_f\}|.$$

The *total demand* is the number of active/pending jobs  $D^t = |\mathcal{D}^t|$ .

2. **Price update.** As discussed in §3.2 below, the market computes  $P^t$  from the load  $\alpha^t$  computed from  $S_f^t$  and  $D^t$ . This updated price  $P^t$  is then posted for the entire period.
3. **Job submission and matching.** Users observe  $P^t$ , decide  $\underline{w}_d(P^t)$ , and enter the matching queue. Jobs are matched to providers via a greedy matching algorithm discussed in §6. Providers are paid according to their reported cost  $\hat{c}_s$  and a premium sharing mechanism given in §4.
4. **Accounting and rollover.** Any unfinished or unassigned jobs remain in the backlog included in  $\mathcal{D}^t$ . Providers' availability  $\tau_s$  decrements unless they may restake.

Once all of these steps have been completed the market proceeds into the next period  $t + 1$ .

### 3.2. Algorithmic Pricing and Market Making

Our pricing mechanism maps a notion of *load* ( $\alpha^t$ ) to a price ( $P^t$ ) at each period  $t$  based on a calculation inspired by Sweeney et al. (2024). When total demand is less than the current floor supply, the pricing mechanism sets the price to  $P_f$ ; otherwise it sets the price to  $P^t = f^t(\alpha^t)$ , where the pricing function  $f^t$  satisfies several standard economic assumptions given below and  $\alpha^t$  is

$$\alpha^t(P^t) := \begin{cases} 1, & D^t \leq S_f^t, \\ \frac{D^t}{S_f^t}, & \text{otherwise.} \end{cases} \quad (1)$$

<sup>2</sup>In §8.1 we discuss how the floor price can be updated rather than just treated as a pre-specified constant.



Intuitively, when demand is at most floor supply ( $\alpha^t = 1$ ) the market can clear at the floor price; otherwise, when demand exceeds the floor supply ( $\alpha^t > 1$ ) the market must raise prices smoothly to activate dormant supply or attract new supply. The reason we measure load using  $S_f^t$  rather than  $S^t = |\mathcal{S}^t|$  is to decouple instantaneous quoting from contemporaneous price-induced supply shifts. Doing so prevents self-reinforcing spikes and simplifies analysis.

For the remainder of the paper, we make the following assumption about the pricing function  $f^t$ .

**Assumption 2.** The pricing function  $f^t$  is continuous, increasing, and concave with floor plateau  $f^t(1) = P_f$  and bounded slope  $\sup_{\alpha \in [1, \infty]} f'^t(\alpha) < \infty$ .

### 3.3. Equilibrium Quotes

Define the cumulative active supply at price  $P^t$  as  $S^t := |\{s \in \mathcal{S}^t : \hat{c}_s \leq P^t\}|$  and be given  $\alpha^t(P^t)$  by equation 1. An *equilibrium quote* at time  $t$  is any solution of

$$P^{t\star} = f^t\left(\alpha^t(P^{t\star})\right). \quad (2)$$

In §5 rigorously explore equilibria in these markets. We prove existence and uniqueness of equilibrium quotes. We further characterize when an equilibrium price  $P^{t\star}$  is *admissible*, i.e.,  $S^t \geq D^t$ , via a “regular crossing” condition on  $D^t - S^t$  near the minimal admissible price and a local responsiveness condition on  $f^t$  at  $\alpha = 1$ .

## 4. Incentive Compatibility

While the algorithmic pricing rule determines an equilibrium market price, this alone does not guarantee that providers reveal their true costs or contribute their full capacity. On the supply side, however, without additional incentives providers may withhold capacity or report inflated costs, leading to inefficiency and higher prices. Thus, we require an explicit *incentive mechanism* to ensure two behaviors: (i) providers stake their *entire available time*, and (ii) they stake at their *lowest acceptable price*. When these behaviors hold, the market operates at full capacity and at the lowest sustainable price. On the demand side, incentives are already aligned: users naturally bid for as much compute as they need at the prevailing price.

To achieve this, our design relies on two key components: *Pool Sharing*, which ensures that providers benefit from staking early and contributing to the collective pool, and *Cheapest Matching*, which guarantees that demand is always allocated first to the lowest-priced available supply. Together, these components align provider incentives with market efficiency, ensuring both truthful cost revelation and maximal supply participation.

### 4.1. Pool-Sharing Scheme

The intuitions behind pool sharing rather than idiosyncratic job matching are twofold: firstly, we want the provider to stake as early as possible rather than wait until the price is high enough to stake. With the pool sharing scheme and our choice of incentive-compatible matching algorithm, we always match the cheapest provider that can accommodate a job’s duration (which we call **CFM-Cheapest Feasible Matching**), providers gain strictly positive payoff by staking as early as possible.

With pool sharing, once a provider  $s$  is matched to a job starting at time  $t$  with duration  $\tau$ , she is compensated hourly from two sources:

1. **Base rate:** her reported cost  $\hat{c}_s$ , paid for each hour worked.



2. **Premium pool:** a share of the surplus generated by all jobs that *start at or after t* and run concurrently during each hour of her work.

Formally, let  $S_{\geq h}^t$  denote the set of providers who are working at hour  $h$  and whose jobs began at or after  $t$  (so providers who started earlier than  $s$  do not share their premiums). Let

$$\Pi_{\geq h}^t = \sum_{s' \in S_{\geq h}^t} (P_{s'} - \hat{c}_{s'})$$

denote the *premium pool* relevant for  $s$  at hour  $h$ , and let  $n_{\geq h}^t = |S_{\geq h}^t|$  be its size.

Then provider  $s$ 's payment at hour  $h$  is

$$\hat{c}_s + \frac{1}{n_{\geq h}^t} \Pi_{\geq h}^t.$$

Her total return for serving  $\tau$  hours is therefore

$$R_s(\hat{c}_s, t, \tau) = \sum_{h=t}^{t+\tau-1} \left( \hat{c}_s + \frac{1}{n_{\geq h}^t} \Pi_{\geq h}^t \right). \quad (3)$$

**Interpretation.** The base rate  $\hat{c}_s$  ensures cost coverage, while the premium share compensates for the opportunity cost of being occupied when later jobs arrive. Importantly,  $s$  does *not* obtain premiums from providers whose jobs started before her match, reflecting that she could not have benefited from those earlier job matching opportunities.

For example, if  $s$  is the only provider in  $S_{\geq h}^t$  ( $n_{\geq h}^t = 1$ ), then  $\Pi_{\geq h}^t = P_s - \hat{c}_s$ , so her hourly payment is  $P_s$ , i.e. the full market price of her job. If there are multiple providers starting at or after  $t$ , each receives her own base rate plus an equal  $1/n_{\geq h}^t$  share of the premium pool  $\Pi_{\geq h}^t$ . This scheme ensures that providers prefer to stake early and at their true cost: early stakers participate in more premium pools, while cheaper providers are always prioritized for matching.

## 4.2. Cheapest Matching

The second central requirement on the supply side is that providers must be incentivized to report their true costs. At a high level, this can be achieved by ensuring that the matching rule *always prioritizes the cheapest eligible provider and true (or subjective belief of provider on) matching competitiveness is high enough* so that the marginal gain from raising reported cost is erased by sharp drop in job matching probability.

The cheapest matching condition can be satisfied by many algorithmic variants, including randomized rules that assign with higher probability to cheaper providers—where again, lowering one's price strictly increases the chance of being matched.

Formally, we make two assumptions for providers' individual rationality and subjective beliefs on matching probability

**Assumption 3** (quasi-rationality). Providers are **quasi-rational** if they never report lower than their true operational cost  $\hat{c}_s < c_s$  because under-reporting will lead to positive probability of negative payoff.

**Assumption 4** (matching competitiveness). Denote provider  $s$ 's subjective belief on the probability that he will be matched with reported cost  $\hat{c}_s$  and staking time  $\tau_s$ ,  $\mathbb{P}_s(s \text{ is matched})|\hat{c}_s, \tau_s$ , as  $\psi_s(\hat{c}_s, \tau_s)$ , such

that the hazard rate of  $\psi_s$ ,  $-\frac{\partial \psi_s(\hat{c}_s, \tau_s)/\partial \hat{c}_s}{\psi_s(\hat{c}_s, \tau_s)}$  is bounded below by  $\frac{\sum_{h=t}^{t+\tau_s-1} (1 - \frac{1}{n_{\geq h}^t})}{R(\hat{c}_s, \tau_s)}$  where  $R_s$  is provider's payoff function from  $t$  to  $t + \tau_s$ .



To explain these two assumptions in hand, we give the proposition that any cheapest matching algorithm is incentive-compatible for provider to report  $\hat{c}_s = c_s$ .

**Proposition 1.** *Any cheapest matching algorithm is incentive-compatible for providers to report  $\hat{c}_s = c_s$  if both Assumption 3 and Assumption 4 are satisfied.*

*Proof.* The provider's payoff (if matched) is given by

$$R_s(\hat{c}, t, \tau) = \sum_{h=t}^{t+\tau-1} \left( \hat{c}_s + \frac{1}{n_{\geq h}^t} \Pi_{\geq h}^t \right).$$

and this profit function is strictly increasing in  $\hat{c}_s$ , where the minimum and maximum are attained at  $c_s$  and  $P$ , respectively.

Since we have proved that under the pool sharing mechanism, it is strictly dominant for the provider to stake all her availability instantly to the network in Proposition 2, therefore  $\tau_s$  is a fixed constant, as any cheapest matching algorithm prioritizes low  $\hat{c}_s$ , we have

$$\frac{\partial \psi_s(\hat{c}_s, \tau_s)}{\partial \hat{c}_s} < 0$$

Take the derivative of provider  $s$  block  $t$ 's expected return of with respect to  $\hat{c}_s$ , we have

$$\frac{\partial \mathbb{E}[R_s(\hat{c}_s, t, \tau_s)]}{\partial \hat{c}_s} = \frac{\partial \psi_s}{\partial \hat{c}_s}(R_s(\hat{c}_s, t, \tau_s)) + \psi_s(\hat{c}_s, \tau_s) \left( \sum_{h=t}^{t+\tau_s-1} \left( 1 - \frac{1}{n_{\geq h}^t} \right) \right)$$

Rearranging, since  $R_s > 0, \psi_s > 0$ , we can get

$$\frac{\partial \mathbb{E}[R_s(\hat{c}, t, \tau_s)]}{\partial \hat{c}_s} = \psi_s \pi(\hat{c}_s, \tau_s) \left( \frac{\psi'_s + \frac{\sum_{h=t}^{t+\tau_s-1} (1 - \frac{1}{n_{\geq h}^t})}{R(\hat{c}_s, \tau_s)}}{\psi_s} \right) \leq 0$$

and the last inequality comes from our second assumption. Since the expected profit is weakly decreasing in  $\hat{c}_s$ , the maximum expected utility is obtained at  $\min \hat{c}_s = c_s$ .  $\square$

We will explain the two assumptions in plain words: (1) Providers never report below their true operational cost. This comes from a quasi-rationality condition, such that they want to eliminate the possibility of negative profit. For example, when provider  $s$ 's  $c_s > P_f$ , but in order to get a match, provider  $s$  reports  $\hat{c}_s = P_f$ , when  $P^t = P_f$ , the pool sharing term is 0, the hourly payoff will equal to  $\hat{c}_s$  which is strictly lower than her true cost. (2) The second condition states that the 'competition intensity' is bounded below at each reported willing price  $\hat{c}_s$ , so that if a provider would raise price, her expected payoff would decrease because of the drop in matching probability, moreover, the smaller the price, the fiercer the competition.

## 5. Equilibrium Quote Price Characterization

In this section, we characterize the equilibrium of our model. Throughout our analysis, we adopt the following assumption regarding the full information of two sides of the market:

**Assumption 5.** We assume that providers are honest i.e. report true tier, and have full information of their submitted jobs. In other words, the provider accurately joins the market tier her machine belongs to and reveals true availability of her time; the user has a correct estimate of her job length and the maximum utility can be attained within her deadline.



This assumption facilitates analytical tractability, allowing us to establish equilibrium existence and derive the optimal matching algorithm. In later sections, we will relax these informational assumptions and examine the implications of incomplete or asymmetric information. We start by formulating the individual optimization problems faced by providers and users, the heterogeneous agent case at  $t$ , and then the existence and uniqueness of the equilibrium quote.

### 5.1. Providers with Availability and Cost

For this section, we will define the single provider's problem and derive its optimal decision, and then show that our pool-sharing mechanism is an incentive compatible mechanism for them to stake as early as they can. In the second subsection, we will derive the distributional effect of optimal decision for heterogeneous providers.

#### 5.1.1. Single Provider

A provider  $s$  present at time  $t$  is characterized by the tuple  $(\tau_s, c_s)$ , where  $\tau_s$  is the availability at  $t$ ,  $c_s$  is the true cost she would provide her compute, and she has to report  $\hat{c}_s$  when joining the market. We assume that providers are speculative; they do not foresee future availability of their compute but aim to maximize the a linear payoff during her entire available time, i.e. the total return from time  $t$  to  $t + \tau_s$ .

**Proposition 2.** *Assume that during provider  $s$ 's  $\tau$  hours of availability, and the provider can stake to the network under two different circumstances, when the market is **over-demanded**, which means instant matching when staked to the market, and when the market is **over-supplied**, which means that the provider may not be matched at the instant of staking. We show that it is:*

1. strictly dominant to stake as early as possible when the market is over-demanded.
2. weakly dominant to stake as early as possible and incentive compatible to stake at  $\hat{c}_s = c_s$  when the market is over-supplied.

*Proof.* When the network is over-demanded, which means that for any  $\hat{c}_s \leq P^t$ , the provider will get matched instantly, then her returns from staking now and  $\tau' < \tau$  periods later are given by

$$R_s(\hat{c}, t, \tau) = \sum_{h=t}^{t+\tau-1} \left( \hat{c}_s + \frac{1}{n_{\geq h}^t} \Pi_{\geq h}^t \right). \quad \text{and} \quad R_s(\hat{c}t + \tau', \tau - \tau') = \sum_{h=t+\tau'}^{t+\tau'-1} \left( \hat{c}_s + \frac{1}{n_{\geq h}^t} \Pi_{\geq h}^t \right).$$

The difference is given by

$$R_s(\hat{c}, t, \tau) - R_s(\hat{c}, t + \tau', \tau - \tau') = \tau' \hat{c}_s + \sum_{h=t}^{t+\tau'-1} \frac{1}{n_{\geq h}^t} \Pi_{\geq h}^t = \sum_{h=t}^{t+\tau'-1} \left( \hat{c}_s + \frac{1}{n_{\geq h}^t} \Pi_{\geq h}^t \right) = R_s(\hat{c}, t, \tau').$$

which is strictly positive, thus it is a strictly dominant strategy to stake earlier than later when the market is over-demanded at any provider's reported willing price  $c_s \leq P^t$ .

When the market is over-supplied, either because the provider's staking price is higher than the market price, i.e.,  $c_s > P^t$ , or the current providers whose reported willing prices lower than provider  $s$  can cover all submitted jobs, which the provider would not be guaranteed with immediate matching once staked, it is still a weakly dominant strategy for the provider to stake earlier rather than later. Since once the market is less over-supplied, either the market price rises to include provider  $s$  automatically, or new jobs are submitted because of the market price  $P$  decreases according to ??, provider  $s$  will be automatically matched once she can be matched for a job.  $\square$



This proposition shows that for whatever  $\hat{c}_s$  the provider reports, it is always optimal for her to join the market as soon as possible, therefore, what is left for us to solve in the matching algorithm is how to incentivize the providers to stake at their true willing price such that  $\hat{c}_s = c_s$ , which we provide a solution in Section 6.

### 5.1.2. Heterogeneous Providers at $t$

In this section, we will show that at any period  $t$  with heterogeneous providers making optimal decision, the number of supplying providers to the market is non-decreasing in the market price  $P$ .

At time  $t$ , let the provider population have distribution  $F^t$  over the reported costs of providers  $\hat{c}_s$ ,

Define the aggregate active provider count at  $P$ :

$$S^t(P) := \int \mathbf{1}\{\hat{c}_s \leq P\} dF^t(\hat{c}),$$

**Proposition 3** (Upward-sloping (weak) aggregate supply for providers). *For any provider's type distribution  $F_s(\theta)$ , if  $P' < P$  then*

$$S^t(P) \geq S^t(P')$$

*Hence, when the per-hour price increases, the number of staking providers does not decrease.*

This is a straightforward result for any distribution  $F_s(\theta)$  as all providers will supply their full availability according to Proposition 2 and  $S^t$  is a non-decreasing function in  $P$ , moreover, the floor supply at  $t$  is given by  $S_f^t = S^t(P_f)$ .

## 5.2. Users with Budget, Deadline, and Discrete Job Hours

At each time  $t$ , the platform posts a per-hour price  $P^t > 0$ . Users decide whether to submit a job and, if so, how many hours to run. Job hours, minimum viable hours, and deadlines are positive integers; jobs must finish *before* the deadline.

### 5.2.1. Single User

A user  $d$  present at time  $t$  is characterized by the tuple  $(B_d, T_d, \underline{w}_d)$ ; where  $B_d > 0$  is the budget,  $T_d \in \mathbb{N}$  is the deadline, and  $\underline{w}_d \in \mathbb{N}$  is the minimum viable run length (some preliminary result from training or partial completion of batch inference). The length of the job is an integer number of hours  $w \in \{0\} \cup \{\underline{w}_d, \underline{w}_d + 1, \dots, \lfloor \frac{B_d}{P^t} \rfloor\}$ , where feasibility requires finishing *before* the deadline and staying within budget.

The user derives value  $v_d : \mathbb{N}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$  from hours, with  $v_d(0) = 0$ . Define discrete marginal gains

$$\Delta_d(w) := v_d(w) - v_d(w-1), \quad w \in \mathbb{N}.$$

**Assumption 6.** We assume two standard properties for the users:

- $v_d$  is increasing and discretely concave;
- If the user is indifferent between submitting and not, she submits (tie-break toward entry).

The single-user optimization at time  $t$  is

$$\max_{w \in \{0\} \cup \{\underline{w}_d, \dots, U_d(P^t)\}} v_d(w) - P^t w. \tag{4}$$

Let  $w_d(P^t)$  be the chosen hours.



**Definition 1.** Define the price-clearing marginal-count

$$h_d(P) := \max \left\{ m \in \mathbb{N}_{\geq 0} : \Delta_d(w) \geq P \text{ for all } w \leq m \right\}$$

which is well-defined and weakly decreasing in  $P$  under Assumption 6. Then

$$w_d(P) = \begin{cases} 0, & \text{if } \lfloor \frac{B_d}{P} \rfloor < \underline{w}_d \text{ or } \max_{w \in \{\underline{w}_d, \dots, U_d(P)\}} [v_d(w) - Pw] \leq 0, \\ \min \left\{ \lfloor \frac{B_d}{P} \rfloor, \max \{\underline{w}_d, h_d(P)\} \right\}, & \text{otherwise,} \end{cases} \quad (5)$$

the price clearing marginal count  $h_d$  can be regarded as the completion of the job, which is when the user's utility is maximized.

**Proposition 4** (Monotone comparative statics, single user). *By Assumption 6. If  $P' < P$ , then*

$$\left\lfloor \frac{B_d}{P'} \right\rfloor \geq \left\lfloor \frac{B_d}{P} \right\rfloor, \quad h_d(P') \geq h_d(P), \quad w_d(P') \geq w_d(P), \quad \mathbf{1}\{w_d(P') > 0\} \geq \mathbf{1}\{w_d(P) > 0\}.$$

Thus, a lower price weakly increases both the probability of submission and the chosen job length.

*Proof.*  $P' < P$  implies  $\lfloor B_d/P' \rfloor \geq \lfloor B_d/P \rfloor$ ; if  $\Delta_d(w) \geq P$ , then  $\Delta_d(w) \geq P'$ , so  $h_d(P') \geq h_d(P)$ ; for any  $w$ ,  $v_d(w) - P'w = (v_d(w) - Pw) + (P - P')w \geq v_d(w) - Pw$ . Together with the previous steps and equation 5, both arguments of min and max weakly increase, implying  $w_d(P') \geq w_d(P)$ ; finally, with the entry tie-break,  $\mathbf{1}\{w_d(\cdot) > 0\}$  is weakly decreasing in  $P$ .  $\square$

### 5.2.2. Heterogeneous Users at $t$

In this section, we will show that at any period  $t$  with heterogeneous users making optimal decisions, the number of jobs submitted to the market is non-increasing in the market price  $P$ .

At time  $t$ , let the user population have distribution  $F^t$  over  $\theta = (B, D, w_0)$  and, optionally, preference shifters inside  $v(\cdot)$  (e.g.,  $v(w) = au(w)$  with  $a > 0$  and discretely concave  $u$ ). Each user solves equation 4 with solution  $w_\theta^*(P^t)$ )

Define the aggregate submitted job count at price  $P$  as,

$$D^t(P) := \int \mathbf{1}\{w_\theta^*(P) > 0\} dF^t(\theta),$$

**Proposition 5** (Downward-sloping (weak) aggregate demand for jobs). *Under Assumption 6 for all types, if  $P' < P$  then*

$$D^t(P') \geq D^t(P) \quad \text{and} \quad H^t(P') \geq H^t(P).$$

Hence, when the per-hour price decreases, neither the number of submitted jobs nor the total compute hours decrease.

*Proof.* By Proposition 4, for each  $\eta$ ,  $\mathbf{1}\{w_\eta^*(P') > 0\} \geq \mathbf{1}\{w_\eta^*(P) > 0\}$  and  $w_\eta^*(P') \geq w_\eta^*(P)$ . Integrating these pointwise inequalities with respect to  $F^t$  preserves the weak inequalities by linearity of the integral. Therefore, no distributional edge case can reverse monotonicity.  $\square$

With discrete hours, hard integer deadlines, and budgets, discrete concavity plus an entry tie-break ensure that lowering price weakly enlarges each user's feasible set and optimal hours, and therefore (weakly) increases both the number of jobs and the aggregate hours in the cross-section.



### 5.3. Equilibrium

Before delving into the equilibrium characterization, we introduce the notations and derive the simple distributional effect of heterogeneous providers and users.

In period  $t$ , denote the set of providers and users as  $\mathcal{S}, \mathcal{D}$ . For each provider  $s \in \mathcal{S}$ , who has a true cost  $c_s$  and availability  $\tau_{t,s}$ , denote the number of providers willing to provide compute at price  $P$  as  $S(P) = |\{s \in \mathcal{S} | \hat{c}_s \leq P\}| \in \mathbb{N}$ . For each user  $d \in \mathcal{D}$ , each will decide whether to submit a job with length  $w_d$ , and let  $D_{t,\min} \geq 0$  denote the number of unfinished jobs at period  $t$ . We first give the definition of an equilibrium quote price of the market.

**Assumption 7** (Environment at time  $t$ ). Fix  $t$  and a floor price  $P_f > 0$ . Let  $b_{\max} \geq P_f$  be a finite upper bound on users' hourly budgets at time  $t$ . Demand and supply are given by functions

$$D^t : [P_f, \infty) \rightarrow (0, \infty), \quad S^t : [P_f, \infty) \rightarrow (0, \infty),$$

that are continuous, with  $D^t$  non-increasing and  $S^t$  non-decreasing in  $P$ , and  $S^t(P_f) > 0$ . (If a backlog of unfinished jobs  $D_{t,\min} \geq 0$  exists, we either include it additively in demand,  $D^t(P) = D_{\text{fresh}}^t(P) + D_{t,\min}$ , or impose  $D^t(P) \geq D_{t,\min}$  for all  $P$ .)

**Definition 2** (Equilibrium quote). Define the load  $\alpha$  as a function of market price

$$\alpha^t(P) = \begin{cases} 1 & \text{if } D^t(P) \leq S^t(P_f) \\ \frac{D^t(P)}{S^t(P_f)} & \text{otherwise} \end{cases} \quad (6)$$

A pricing rule is any function  $f^t : [1, \infty) \rightarrow [P_f, b_{\max}]$  that is continuous and strictly increasing, and satisfies  $f^t(1) \geq P_f$ . An *equilibrium quote* at time  $t$  is any  $P^{t\star} \geq P_f$  solving the fixed-point equation

$$P^{t\star} = f^t(\alpha^t(P^{t\star}))$$

**Proposition 6** (Existence and uniqueness of equilibrium quote). *Under the assumptions above, for every pricing rule  $f^t$  as in the definition, there exists a unique equilibrium quote  $P^{t\star} \in [P_f, b_{\max}]$ . Moreover, the load  $\alpha^t(P)$  is non-increasing in  $P$ .*

*Proof.* Let  $\Phi(P) := P - f^t(D^t(P)/S^t(P_f))$  on  $[P_f, \infty)$ . By continuity of  $D^t$  and  $f^t$ ,  $\Phi$  is continuous. Since  $D^t$  is non-increasing and  $f^t$  is strictly increasing,  $f^t(D^t(P)/S^t(P_f))$  is non-increasing in  $P$ ; hence  $\Phi$  is strictly increasing. At  $P_f$ ,

$$\Phi(P_f) = P_f - f^t\left(\frac{D^t(P_f)}{S^t(P_f)}\right) \leq P_f - f^t(1) \leq 0,$$

and for any  $\tilde{P} \geq b_{\max}$ ,

$$\Phi(\tilde{P}) = \tilde{P} - f^t\left(\frac{D^t(\tilde{P})}{S^t(P_f)}\right) \geq \tilde{P} - b_{\max} \geq 0.$$

By the intermediate value theorem, there exists  $P^{t\star} \in [P_f, b_{\max}]$  with  $\Phi(P^{t\star}) = 0$ . Uniqueness follows from strict monotonicity of  $\Phi$ . Finally,  $\alpha^t(P)$  is non-increasing in  $P$  because  $S^t(P_f)$  is constant and  $D^t$  is non-increasing.  $\square$

**Definition 3** (Admissibility threshold). The (lowest) admissible price level is

$$P_{\text{adm}}^t := \inf\{P \geq P_f : S^t(P) \geq D^t(P)\},$$

assumed finite. A price  $P$  is *admissible* if  $S^t(P) \geq D^t(P)$ .



**Proposition 7** (Admissibility of the equilibrium quote). *Fix  $t$  and let  $P_f > 0$  with  $S^t(P_f) > 0$ , and let the pricing rule  $f^t : [1, \infty) \rightarrow [P_f, b_{\max}]$  be continuous, non-decreasing, with a floor plateau  $f^t(\alpha) = P_f$  for  $\alpha = 1$ , strictly increasing on  $[1, \infty)$ , and normalized by  $f^t(1) = P_f$ . Let*

$$P_{\text{adm}}^t := \inf\{P \geq P_f : S^t(P) \geq D^t(P)\}.$$

Assume:

- (S1) (Regular crossing) There exists  $\delta_t > 0$  such that for all  $P \in [P_f, P_{\text{adm}}^t]$ ,  $D^t(P) - S^t(P) \geq \delta_t (P_{\text{adm}}^t - P)$ .  
(F1) (Local responsiveness) The right derivative  $(f^t)'(1^+)$  exists and  $(f^t)'(1^+) \geq S^t(P_f)/\delta_t$ .

Then the unique fixed point  $P^{t*}$  solving  $P^{t*} = f^t(\alpha^t(P^{t*}))$  is admissible:  $P^{t*} \geq P_{\text{adm}}^t$  (hence  $S^t(P^{t*}) \geq D^t(P^{t*})$ ).

*Proof.* Write  $D := D^t$ ,  $S := S^t$ ,  $S_f := S^t(P_f)$ ,  $\alpha(P) := D(P)/S_f$ ,  $f := f^t$ , and  $P_{\text{adm}} := P_{\text{adm}}^t$ . For  $P \leq P_{\text{adm}}$ , (S1) and  $S(P) \geq S_f$  imply

$$\alpha(P) = \frac{D(P)}{S_f} \geq \frac{S(P) + \delta_t(P_{\text{adm}} - P)}{S_f} \geq 1 + \frac{\delta_t}{S_f} (P_{\text{adm}} - P).$$

By (F1) and the mean value inequality at  $\alpha = 1$ ,

$$f(\alpha(P)) \geq f(1) + (f'(1^+))(\alpha(P) - 1) \geq P_f + \frac{S_f}{\delta_t} \cdot \frac{\delta_t}{S_f} (P_{\text{adm}} - P) = P_f + (P_{\text{adm}} - P).$$

Since any feasible price satisfies  $P \geq P_f$ , we obtain for all  $P \in [P_f, P_{\text{adm}}]$ :

$$f(\alpha(P)) \geq P_f + (P_{\text{adm}} - P) \geq P_{\text{adm}} > P.$$

Hence no fixed point can lie below  $P_{\text{adm}}$ . Therefore the unique fixed point  $P^*$  satisfies  $P^* \geq P_{\text{adm}}$ , i.e., it is admissible.  $\square$

With existence and uniqueness of the time-wise equilibrium quote  $P^{t*}$  secured, the remaining design problem is normative: choose a floor  $P_f$  such that the endogenous fixed point typically resides at, or arbitrarily close to, the floor. We select  $P_f$  so that it adheres to three criteria: (i) keep the utilization ratio  $\alpha^t = D^t/S_f$  at or below one with high probability, thereby suppressing price excursions; (ii) satisfy participation constraints for a thick mass of low-cost providers; and (iii) allow users to lock in full jobs at the posted floor. We formalize what it means to be a "healthy" floor price in Section 8.1.

## 6. Welfare & Optimal Matching Design

As we pinned down the equilibrium quote, we will proceed to define the welfare of the market. For providers, they aim to maximize the monetary value from renting their available compute time, and for users, their aggregate welfare is given by the proportion of completed jobs upon matching.

At  $t$ , the market price is given by  $P^t$ , we denote the available providers by  $s \in \mathcal{S}^t$ , where each of them can be characterized by two characteristics, available time  $\tau_s$  and reported cost  $\hat{c}_s$ . From user (job) side, we denote the submitted job as  $d \in \mathcal{D}^t$ , which is a queue of jobs characterized by purchasing orders with time  $w_d^t$ . The matching  $(s, d) \in \mathcal{M}^t$  is a partial bijection:

$$\mathcal{M}^t \subset \mathcal{S}^t \times \mathcal{D}^t.$$

Each provider can be matched at most once, since we are focusing on the single period optimal matching first, we will omit  $t$  for brevity until the multi-period analysis.

One important definition of this matching market is feasibility of a provider  $s$  to job  $d$ , formally defined as:



**Definition 4.** A provider  $s$  with staking time  $\tau_s$  is feasible for a certain job with length  $w_d$  if and only if

$$\tau_s \geq w_d$$

which indicate that the matched provider  $s$  can finish the job. Since whether a matched job is feasible affects the welfare of users, we use this metric to define the welfare of a user for a given matching  $\mathcal{M}$ .

**Definition 5.** The welfare consists of two optimization objectives, first, minimizing cost objective, which is the accumulated reported cost of the matched provider under matching  $\mathcal{M}$ :

$$\min \sum_{(s,d) \in \mathcal{M}} \hat{c}_s \quad (7)$$

and second, the completion objective, which is the number of feasible matched jobs under matching  $\mathcal{M}$ :

$$\max \sum_{(s,d) \in \mathcal{M}} \mathbf{1}_{\{w_d - \tau_s \leq 0\}}. \quad (8)$$

For any matching algorithm  $\mathcal{M}$ , we define **S(supply)**, **D(demand)**-regret as the welfare loss compared to the optimal **online** matching that minimizes Equation 8 and Equation 7. Denote the matching produced by optimal matching algorithm that minimizes Equation 8 and Equation 7 as  $\mathcal{M}_{S-opt}$ ,  $\mathcal{M}_{D-opt}$ , respectively, we have

$$R_S(\mathcal{M}) = \sum_{(s,d) \in \mathcal{M}} \hat{c}_s - \sum_{(s,d) \in \mathcal{M}_{S-opt}} \hat{c}_s \quad (9)$$

and

$$R_D(\mathcal{M}) = - \left( \sum_{(s,d) \in \mathcal{M}} \mathbf{1}_{\{w_d - \tau_s \leq 0\}} - \sum_{(s,d) \in \mathcal{M}_{D-opt}} \mathbf{1}_{\{w_d - \tau_s \leq 0\}} \right). \quad (10)$$

We model the problem as an *online matching* problem following the assumption of computation job scheduling literature (Bao et al., 2018; Zhang et al., 2025; Deng et al., 2021). More specifically, the decentralized nature of our market falls under the category of *online bipartite matching*, which was first introduced by Karp et al. (1990). For online bipartite matching problem, there is a bipartite graph  $G(S, \mathcal{D}, E)$  where one side  $S$  is known to us in advance (the providers) and the other side  $\mathcal{D}$  (the jobs) arrives online, one vertex at a time, when a vertex  $d \in \mathcal{D}$  arrives, its neighbors in  $S$  are revealed. The arriving vertex can be matched to some available neighbor (if any). A match, once made, cannot be revoked. The objective is to maximize the size of the matching. For our market we assume that the arrival of jobs follows an adversarial model as studied by Karp et al. (1990), that is, the vertices (jobs) in  $\mathcal{D}$  arrive online in a preselected order; or a Random Order Model, where the orders arrive online in a random order, that is, there is a uniformly random permutation  $\sigma$  of  $\mathcal{D}$ . For the rest of the section, we stick to the adversarial arrival model and we leave the random order model as an open question for future work.

## 6.1. Provider Welfare Optimization

To optimize the provider-side objective Equation 7, we first show an equivalence theorem stating that for a matched provider  $s$ , minimizing the total cost of other matched providers  $s' \neq s$  is equivalent to maximizing her payoff in the same period.

**Lemma 1.** Denote the total number of matched providers by  $n$ , and price by  $P$ , we omit  $t$  for brevity, thus minimizing

$$\sum_{(s' \neq s, d) \in \mathcal{M}} \hat{c}_{s'}$$



is equivalent to maximizing the time  $t$  payment to a working provider  $s$  if  $s$  is already matched

$$\hat{c}_s + \frac{1}{n} \sum_{(s' \neq s, d) \in \mathcal{M}} (P - \hat{c}_{s'})$$

*Proof.* For each provider  $s$ ,  $\hat{c}_s$ ,  $P$  are fixed in period  $t$  and the second term in the period  $t$  payment can be rewritten as

$$\frac{1}{n} (n \cdot P - \sum_{(s' \neq s, d) \in \mathcal{M}} \hat{c}_{s'})$$

since  $P^t$  and  $n$  are fixed at time  $t$ , then minimizing  $\sum_{(s, d) \in \mathcal{M}} \hat{c}_s$  is equivalent to maximizing provider  $s$ 's time  $t$  payment. □

Building on the general bipartite matching problem, one strict generalization called **Online Vertex-Weighted Bipartite Matching**, which firstly studied by Aggarwal et al. (2011), fits better with our scenario, and the problem can be described as follows: Each vertex  $s \in S$  has a non-negative weight which is known in advance, and the goal is to maximize the sum of weights of vertices in  $S$  that get matched. In our scenario, the minimization problem is an online vertex bipartite matching with (1) a fully connected graph and (2) the weight of each vertex  $s \in S$  being the negative of the cost. Thus, the optimal matching algorithm becomes straightforward – we only have to greedily match the cheapest available provider and this achieves the maximum sum of weights (minimum sum of costs). We will outline the matching algorithm below and refer to it as Greedy-Cheapest Matching(**GCM**).

**Greedy-Cheapest Matching** The Greedy-Cheapest Matching(**GCM**) is given as

- Choose a permutation  $\sigma : \{1, \dots, m\} \rightarrow S$  such that

$$\hat{c}_{\sigma(1)} \leq \hat{c}_{\sigma(2)} \leq \dots \leq \hat{c}_{\sigma(m)}$$

(ascending willing price; ties arbitrary).

- For every job index  $d$ ,  $S_d = S \setminus \text{dom}(\mathcal{M})$ , such that if  $S_d \neq \emptyset$ , select

$$s_i^*(d) := \operatorname{argmin}_{s \in S_d} \hat{c}_s \quad (\text{break ties by the position of } s \text{ in } \sigma)$$

Add the pair to the matching:

$$\mathcal{M} \leftarrow \mathcal{M} \cup \{(s_i^*(d), d)\}.$$

In short, **GCM** always matches the cheapest provider available to the arriving job irrespective of the feasibility, this way, the payoff to the matched providers are maximized.

## 6.2. User Welfare Optimization

As for the optimization of users' (jobs') objective function Equation 8, the problem becomes a little bit more complicated; since for all  $(d, s) \in \mathcal{M}$ , to achieve as many  $w_d - \tau_s \geq 0$  pairs as possible, we have to maximize the number of matched providers whose time staking  $\tau_s$  is larger than or equal to the length of matched job. This problem corresponds to **Online Bipartite Matching of Convex Graph**, we first give the definition of a convex graph  $G(S, \mathcal{D}, E)$  and then elaborate on how user welfare maximization is equivalent to the maximal matching.



**Definition 6.** Let  $G(\mathcal{S}, \mathcal{D}, E)$  be a bipartite graph with partite sets  $\mathcal{S}$  and  $\mathcal{D}$  and edge set  $E \subseteq \mathcal{S} \times \mathcal{D}$ . If there exists a permutation  $\sigma : \{1, \dots, m\} \rightarrow \mathcal{S}$  such that  $s \in \mathcal{S}$  are ordered as

$$s_1 \prec s_2 \prec \dots \prec s_m.$$

For  $d \in \mathcal{D}$ , the neighborhood of  $d$  is defined as

$$N_{\mathcal{S}}(d) := \{s \in \mathcal{S} : (s, d) \in E\}.$$

The graph  $G(\mathcal{S}, \mathcal{D}, E)$  is convex on  $\mathcal{S}$  with respect to the permutation  $s_1 \prec \dots \prec s_m$  if and only if for every job  $d \in \mathcal{D}$ , either  $N_{\mathcal{S}}(d) = \emptyset$  or there exist integers  $a(d), b(d)$  satisfying  $1 \leq a(d) \leq b(d) \leq m$  such that

$$N_{\mathcal{S}}(d) = \{s_{a(d)}, s_{a(d)+1}, \dots, s_{b(d)}\}$$

With this definition, we can draw an analogy of our matching problem. Using the Definition 4 of feasibility, we have the following lemma.

**Lemma 2.** Order providers by nondecreasing  $\tau$  and let  $E = \{(s, d) \in \mathcal{S} \times \mathcal{D} : \tau_s \geq w_d\}$ . Then  $G(\mathcal{S}, \mathcal{D}, E)$  is convex on  $\mathcal{S}$  and for every  $d$ ,  $N_{\mathcal{S}}(d) = \{s_{a(d)}, \dots, s_m\}$  (so  $b(d) = m$ ). Moreover, the following are equivalent:

1. maximizing the number of feasible pairs  $(s, d)$  with  $\tau_s \geq w_d$  in a matching;
2. computing a maximum-cardinality matching in  $G(\mathcal{S}, \mathcal{D}, E)$ .

*Proof.* Convexity and  $b(d) = m$  are immediate from the definition of  $E$ : for fixed  $d$ , if  $s_j$  can process  $d$ , then so can any  $s_{j'}$  with  $j' > j$ , hence  $N_{\mathcal{S}}(d)$  is a suffix interval.

For a matching  $\mathcal{M}$ , let  $\mathcal{M}_{\text{feas}} = \mathcal{M} \cap E$  and  $\mathcal{M}_{\text{infeas}} = \mathcal{M} \setminus E$ . Then  $|\mathcal{M}| = |\mathcal{M}_{\text{feas}}| + |\mathcal{M}_{\text{infeas}}|$ . Among such matchings, maximizing  $|\mathcal{M}_{\text{feas}}|$  is equivalent to maximizing  $|\mathcal{M}|$  subject to  $\mathcal{M} \subseteq E$ , which is exactly (2).  $\square$

According to [Glover \(1967\)](#), the following matching algorithm attains the maximal-cardinality matching for the convex bipartite graph  $G(\mathcal{S}, \mathcal{D}, E)$  thus equivalently maximizes the objective function 8 according to Lemma 2.

**Greedy-Shortest Matching** The Greedy-Shortest Matching (**GSM**) is given as

- Choose a permutation  $\sigma : \{1, \dots, m\} \rightarrow \mathcal{S}$  such that

$$\tau_{\sigma(1)} \leq \tau_{\sigma(2)} \leq \dots \leq \tau_{\sigma(m)}$$

(ascending staking time; ties arbitrary).

- For every job index  $j$ ,  $\mathcal{S}_j = \mathcal{S} \setminus \text{dom}(\mathcal{M})$ , such that if  $\mathcal{S}_j \neq \emptyset$ , select

$$s_i^*(d) := \operatorname{argmin}_{s \in \mathcal{S}_j} \tau_{s_i} \quad \text{s.t.} \quad \tau_{s_i} \geq w_j$$

Add the pair to the matching:

$$\mathcal{M} \leftarrow \mathcal{M} \cup \{(s_i^*(d), d)\}.$$

Otherwise, reject  $d$ .

In short, we always match the shortest-feasible provider available to the current job, then we show the maximal cardinality of feasible matching of **GSM**.



**Proposition 8.** Let  $\mathcal{M}_{GSM}$  be the final matching produced by **GSM** and  $n := |\mathcal{D}|$ . Define  $t(d) = \min \{i : \tau_{s_i} \geq w_d\}$  if it exists; else  $t(d) = +\infty$ . Feasible providers for  $d$  are  $N_S(d) = \{s_i : i \geq t(d)\}$ . For  $s = 1, \dots, m+1$  define the suffix sets

$$\mathcal{S}_{\geq i} := \{s_i, \dots, s_m\} \quad (\mathcal{S}_{\geq m+1} = \emptyset), \quad \mathcal{D}_{\geq i} := \{d \in \mathcal{D} : t(d) \geq i\},$$

and the deficiency  $\Delta$

$$\Delta := \max_{1 \leq s \leq m+1} (|\mathcal{D}_{\geq i}| - |\mathcal{S}_{\geq i}|)_+.$$

Then we have,

1.  $\mathcal{M}_{GSM}$  has the maximum cardinality among all matchings of  $G(\mathcal{S}, \mathcal{D}, E)$ .
2.  $|\mathcal{M}_{GSM}| = n - \Delta$ , where equivalently

$$\frac{|\mathcal{M}_{GSM}|}{n} = 1 - \frac{\Delta}{n} = 1 - \frac{1}{n} \max_{1 \leq s \leq m+1} (|\mathcal{D}_{\geq i}| - |\mathcal{S}_{\geq i}|)_+.$$

Firstly we prove the following lemma

**Lemma 3.** If **GSM** rejects a job in  $\mathcal{D}_{\geq i}$ , then every provider in  $\mathcal{S}_{\geq i}$  is matched in the final **GSM** matching.

*Proof.* Let  $d \in \mathcal{D}_{\geq i}$  be the first job (in time) that **GSM** rejects among those with threshold  $\geq i$ . At its arrival all providers in  $\mathcal{S}_{\geq i}$  are already taken-otherwise **GSM** would have matched  $d$  to the available provider with shortest  $\tau_s$ . Since providers, once matched, never become free, by the end the whole  $\mathcal{S}_{\geq i}$  is matched.  $\square$

*Proof for Proposition 8.* Using Lemma 3, we show that **GSM** produces a matching  $\mathcal{M}_{GSM}$  with  $|\mathcal{M}_{GSM}| = n - \Delta$ .

Process the arrival sequence  $d^1, \dots, d^{(t)}$ . For  $t = 0, 1, \dots, n$ ,  $(t)$  here represents the arrival order of jobs within one period, to avoid confusion, we use  $(t)$ . We maintain a maximum matching  $\mathcal{M}^{(t)}$  such that

$$\mathcal{M}^{(t)} \cap \left( \{d_1, \dots, d^{(t)}\} \times \mathcal{S} \right) = \mathcal{M}_{GSM} \cap \left( \{d_1, \dots, d_n\} \times \mathcal{S} \right),$$

and, among all maximum matchings with this property,  $\mathcal{M}^{(t)}$  minimizes lexicographically the multiset of provider indices used on the remaining jobs  $\{d_{t+1}, \dots, d_n\}$ , that is to say, pick  $\mathcal{M}^{(t)}$  that uses the smallest provider indices first for the jobs that are yet to come.

The basis  $t = 0$  holds (pick any maximum matching). Assume  $\mathcal{M}^{(t-1)}$  is fixed; we construct  $\mathcal{M}^{(t)}$  depending on **GSM**'s action on  $d^{(t)}$ .

**Case A (GSM rejects  $d^{(t)}$ ).** By Lemma 3 with  $i = t(d^{(t)})$ , **GSM** has matched all providers in  $\mathcal{S}_{\geq t(d^{(t)})}$ . If  $\mathcal{M}^{(t-1)}$  matched  $d^{(t)}$ , it would need some provider in  $\mathcal{S}_{\geq t(d^{(t)})}$ ; but those are already used by the first  $t-1$  jobs in  $\mathcal{M}_{GSM}$ , hence (by the induction invariant) also in  $\mathcal{M}^{(t-1)}$ . Restricting to  $\mathcal{D}_{\geq t(d^{(t)})} \cup \mathcal{S}_{\geq t(d^{(t)})}$  this would contradict the maximality of  $\mathcal{M}^{(t-1)}$ . Thus no maximum matching agreeing with the prefix can match  $d^{(t)}$ ; set  $\mathcal{M}^{(t)} := \mathcal{M}^{(t-1)}$ .

**Case B (GSM matches  $d^{(t)}$  to  $s_{i^*}$ , the leftmost free feasible provider).** If  $\mathcal{M}^{(t-1)}$  already matches  $d^{(t)}$  to  $s_{i^*}$ , set  $\mathcal{M}^{(t)} := \mathcal{M}^{(t-1)}$ . Otherwise, we build an alternating walk and flip it.

- Start with the non-  $\mathcal{M}^{(t-1)}$  edge  $(s_{i^*}, d^{(t)})$ .
- If  $s_{i^*}$  is unmatched in  $\mathcal{M}^{(t-1)}$ , we will shortly flip a path ending at  $s_{i^*}$ . Otherwise, let  $d^{(1)}$  be the job matched to  $s_{i^*}$  in  $\mathcal{M}^{(t-1)}$ ; note that  $d^{(1)}$  arrives after  $d^{(t)}$  because  $s_{i^*}$  is unused by the prefix in both  $\mathcal{M}_{GSM}$  and  $\mathcal{M}^{(t-1)}$ .



- Follow the  $\mathcal{M}_{GSM}$ -edge incident to  $d^{(1)}$ , which (by **GSM**) is  $(s_{i_1}, d^{(1)})$  where  $s_{i_1}$  is the leftmost feasible provider for  $d^{(1)}$ .
  - If  $s_{i_1}$  is unmatched in  $\mathcal{M}^{(t-1)}$ , stop; else continue: from  $s_{i_1}$  take the  $\mathcal{M}^{(t-1)}$ -edge to some  $d^{(2)}$ , then the  $\mathcal{M}_{GSM}$ -edge to  $s_{i_2}$ .
- If every job is feasible and  $|\mathcal{D}_{\geq i}| \leq |\mathcal{S}_{\geq i}|$  for all  $s$ , then  $\Delta = 0$  and **GSM** matches all jobs

Because each job's neighborhood is a suffix and  $\mathcal{M}_{GSM}$  always uses the leftmost feasible provider for that job, whenever we move at some step from a  $\mathcal{M}_{GSM}$ -edge  $(s_{i_q}, d^{(q)})$  to the  $\mathcal{M}^{(t-1)}$ -edge incident to  $d^{(q)}$ , the provider index cannot decrease:

$$s^* \leq s_1 \leq s_2 \leq \dots$$

The sequence of indices is bounded by  $m$ . If the walk did not reach an  $\mathcal{M}^{(t-1)}$ -free provider, it must eventually repeat a provider index, forming an alternating cycle whose jobs all arrive after  $d^{(t)}$ . Flipping that cycle yields another maximum matching that still agrees with the prefix but uses a lexicographically smaller multiset of provider indices on the remaining jobs (since along the cycle we replace some  $\mathcal{M}^{(t-1)}$ -edges at larger indices by  $\mathcal{M}_{GSM}$ -edges at no larger indices, with at least one strict decrease coming from  $s_{i^*}$ ). This contradicts the choice of  $\mathcal{M}^{(t-1)}$ . Hence the walk reaches an  $\mathcal{M}^{(t-1)}$  free provider  $s_{i_T}$ .

Flip the alternating path from  $d^{(t)}$  to  $s_{i_T}$ . The result is a maximum matching  $\mathcal{M}^{(t)}$  that still matches the prefix exactly as  $\mathcal{M}_{GSM}$  does and now matches  $d^{(t)}$  to  $s_{i^*}$ . This completes the induction. For  $t = n$ , we obtain a maximum matching  $\mathcal{M}^{(n)}$  that coincides with  $\mathcal{M}_{GSM}$ ; therefore  $\mathcal{M}_{GSM}$  is a maximum-cardinality matching.

Therefore, every matching has size at most  $n - \Delta$  and the proposition shows **GSM** achieves this bound, so  $|\mathcal{M}_{GSM}| = n - \Delta$ .  $\square$

We have proved that when the objective is to maximize the number of feasible matchings, **GSM** is the optimal in the adversarial model. However, the problem with **GSM** is that it is not incentivizing providers to stake at their true cost even under Assumption 3 and 4. Therefore, we propose an algorithm that is incentive-compatible and achieves bounded D- and S-regret with respect to the **GCM** and **GSM**, respectively, which we call **Cheapest Feasible Matching(CFM)**.

### 6.3. A Incentive-Compatible, Feasibility-Aware Matching Algorithm

**Cheapest-Feasible Matching** The Cheapest-Feasible Matching(**CFM**) is given as

- Choose a permutation  $\sigma : \{1, \dots, m\} \rightarrow \mathcal{S}$  such that

$$\tau_{\sigma(1)} \leq \tau_{\sigma(2)} \leq \dots \leq \tau_{\sigma(m)}$$

(ascending availability; ties arbitrary).

- For every job index  $d = 1, 2, \dots, n$  (queue order): Define the set of still-unused, feasible providers

$$\mathcal{F}_j := \{s \in \mathcal{S} \setminus \text{dom}(\mathcal{M}) : \tau_s \geq w_j\}$$

such that

$$\text{dom}(\mathcal{M}) := \{s \in \mathcal{S} : \exists d \in \mathcal{D}(s, d) \in \mathcal{M}\}$$

is the domain of the relation. If  $\mathcal{F}_j \neq \emptyset$ , select

$$s_i^*(d) := \operatorname{argmin}_{s \in \mathcal{F}_j} \hat{c}_s \quad (\text{break ties by the position of } s \text{ in } \sigma).$$



Add the pair to the matching:

$$\mathcal{M} \leftarrow \mathcal{M} \cup \{(s_i^*(d), d)\}.$$

If  $\mathcal{F}_j = \emptyset$ ,

$$s_i^*(d) := \operatorname{argmax}_{s \in \mathcal{S}} \tau_s \quad (\text{break ties by lower } \hat{c}_s).$$

In short, we sort the available providers by time staking  $\tau_s$ , for every job  $d$ , and select the cheapest provider that can cover their entire job; if not, then select the longest possible provider for job  $d$ . The detailed breakdown of the computational complexity for each algorithm is given in Appendix 10.1.

## 6.4. Regret Analysis

### 6.4.1. Single-Period Regret

**CFM v.s. GCM** To evaluate how our **CFM** algorithm compares to the two optimal algorithms in minimizing cost and unfinished jobs, we dedicate this section to the regret analysis by comparing **CFM** to **GCM** and **GSM**.

We give our first proposition on the regret bound between **CFM** and **GCM**:

**Proposition 9.** *For any job arriving sequence at period  $t$ , and the number of available providers and jobs as  $m, n$ , respectively. For each period, separate the characterization by market price  $P^t$ .*

- When  $P = P_f$ ,

$$\operatorname{cost}_{\text{CFM}} - \operatorname{cost}_{\text{GCM}} \leq 0$$

- When  $P^t > P_f$ ,

- if  $n \leq \lfloor m/2 \rfloor$ , then the S-regret is bounded by  $\operatorname{cost}_{\text{CFM}} - \operatorname{cost}_{\text{GCM}} \leq (P - P_f)n$ .
- if  $m > n > \lfloor m/2 \rfloor$ , then the S-regret is bounded by  $\operatorname{cost}_{\text{CFM}} - \operatorname{cost}_{\text{GCM}} \leq (P - P_f)(m - n)$ .
- if  $n \geq m$ , then the bound goes to 0 because all providers will be matched.

Therefore when  $P^t > P_f$ , the universal bound is given by  $\lfloor \frac{m}{2} \rfloor (P - P_f)$

*Proof.* Consider any job sequence at time  $t$ , for any job  $d$ .

When  $P = P_f$ , then all matched providers have  $\hat{c}_s = P_f$ , and **CFM** matches feasible ones while **GCM** randomizes provider-job matching, therefore **CFM** weakly dominates **GCM**.

When  $P > P_f$ , **GCM** pays at least  $P_f$ ; **GCM** pays the cheapest among feasible, whose cost is  $\leq P^t$ . So

$$0 \leq \operatorname{cost}_{\text{CFM}}(d) - \operatorname{cost}_{\text{GCM}}(d) \leq P - P_f$$

In the worst-case scenario, the shortest job is infeasible for any available provider, and **CFM** matches all providers with  $\hat{c}_s = P^t$ , and for **GCM** the same set of jobs are matched with providers registered at  $\hat{c}_s = P_f$  by **GCM**. However, this only happens when the number of jobs is less than or equal to  $\lfloor \frac{m}{2} \rfloor$ . When the number of jobs is greater than  $\lfloor \frac{m}{2} \rfloor$ , **CFM** and **GCM** will inevitably match at least  $m - (2(m - n)) = 2n - m$  identical providers. In the case that **CFM** matches the most expensive ones and **GCM** matches the cheapest ones, and the remaining  $n - (2n - m) = m - n$  matched providers will differ at most by  $P - P_f$ , thus the cost is bounded above by  $(P - P_f)(m - n)$  in this case, and this bound is maximized at  $n = \lfloor \frac{m}{2} \rfloor$ .  $\square$

**CFM v.s. GSM** Formally, the D-regret between **CFM** and **GSM** is given by

$$R_D(\mathcal{M}_{\text{CFM}}) = - \left( \sum_{(s,d) \in \mathcal{M}_{\text{CFM}}} \mathbf{1}_{\{w_d - \tau_s \leq 0\}} - \sum_{(s,d) \in \mathcal{M}_{\text{GSM}}} \mathbf{1}_{\{w_d - \tau_s \leq 0\}} \right)$$



Firstly, we give an equivalence condition on the provider's reported price and staking time distribution  $(\hat{c}, \tau)$  when those **CFM** and **GSM** coincide. Also from then on, we assume providers are indexed by nondecreasing  $\tau$ .

**Lemma 4** (Coincidence under monotone costs). *Assume providers are indexed by nondecreasing  $\tau$  and **CFM** breaks ties to the left among equal-cost feasible providers. If  $c_1 \leq c_2 \leq \dots \leq c_m$ , then **CFM** and **GSM** produce identical matchings on every arrival sequence; the **D**-regret is zero.*

*Proof.* For any job  $d$ , among feasible providers  $\{s \geq t(d)\}$  the minimum cost is attained at the leftmost feasible index  $t(d)$  because  $c_1 \leq \dots \leq c_m$ . With left tie-breaking, **CFM** chooses  $t(d)$ , which is exactly **GSM**'s choice.  $\square$

This raises the question of when do **CFM** and **GSM** disagree? Intuitively, when a provider with a longer staking time registers a lower cost compared to some provider with a shorter staking time, at some point in time, **GSM** will assign an incoming job to the shorter provider which is more expensive, but **CFM** will assign the same incoming job to a longer but cheaper provider, and this is when **D**-regret accumulates. With this in mind, we define the worst-case provider distribution for **CFM** regarding **D**-regret which we call *anti-sorted* provider distribution.

**Definition 7.** If for every pair of providers  $s_i, s_j$  with initial staking time and cost  $\tau_{s_i}(0), \tau_{s_j}(0), c_{s_i}, c_{s_j}$ , if  $\tau_s \geq \tau_{s'}$ , then  $c_s \leq c_{s'}$ . If the providers' index are sorted by staking time, then for any  $s' < s$ ,  $c_{s'} > c_s$ , that is, for any  $s \in \{1, \dots, m\}$ , it acts as a threshold such that for all  $s' \geq i$ ,  $c_{s'} \leq c_s$ , with this, we call the provider's distribution as *anti-sorted*.

This definition shows that if the index is sorted by time, then every provider can serve as a "cut" such that any providers with a smaller index are "short and expensive," and any provider with a larger index is "longer and cheaper". The anti-sorted provider distribution is one necessary condition for the worst case **D**-regret.

**Theorem 1 (1/2 Competitive Ratio).** *Assume providers are indexed by nondecreasing staking time  $\tau_1 \leq \dots \leq \tau_m$ , each can start at most one job in the period, and at most  $m$  jobs arrive online in an adversarial order. Let **GSM** match each job to the leftmost feasible provider and let **CFM** match each job to a feasible provider of minimum reported cost, breaking ties to the left. Suppose  $n \leq m$  jobs arrive in an online manner, then **D**-regret in the period satisfies*

$$R_D(\mathcal{M}_{CFM}) \leq \left\lfloor \frac{n}{2} \right\rfloor.$$

*Proof.* For every incoming job  $d$ , except when **GSM** and **CFM** agree, there are two types of events that **GSM** and **CFM** will match job  $d$  differently. For a job  $d$  with threshold index  $t(d)$ , exactly one of the following events may occur:

- (Skip) **CFM skips left:** **GSM** matches  $d$  to index  $i = t(d)$ , while **CFM** matches  $d$  to some  $k > i$ . We call  $d$  a *skip job* and record the *skip pair*  $(i, k)$  produced by  $d$ .
- (Rej) **CFM-only rejection:** **GSM** matches  $d$  but **CFM** rejects  $d$ . We call  $d$  a *rejection job*.

The key fact is that every rejection needs an earlier skip that covers it. Assume a rejection occurs on a job  $d^*$  with threshold  $r := t(d^*)$ . Just before processing  $d^*$ , **CFM** has exhausted all providers in the suffix  $\{r, r+1, \dots, m\}$  (otherwise **CFM** could match  $d^*$ ), while **GSM** has not exhausted that suffix (since **GSM** matches  $d^*$ ). Hence, among providers in  $\{r, \dots, m\}$ , **CFM** has used *strictly more* than **GSM**.



Therefore, there exists a previous job  $d$  for which **CFM** used some  $k \geq r$  while **GSM** used some  $i < r$  (otherwise the counts in the suffix would be the same). For that job  $d$ , **GSM** chose its leftmost feasible index  $i = t(d)$  and **CFM** chose  $k > i$ , i.e.,  $d$  is a *skip job* that produced the skip pair  $(i, k)$  with

$$i < r \leq k.$$

We say that this skip *covers* the rejection at threshold  $r$ .

Process the jobs in arrival order. Maintain a set of *unused skip jobs* or, equivalently, one *credit* per skip job. When a rejection occurs at threshold  $r$ , select the *latest* unused skip job whose skip pair  $(i, k)$  covers  $r$  (as shown above, at least one such skip exists), and *charge* the rejection to that skip; then mark the skip as *used*. This produces an *injective* map from rejection jobs to skip jobs because each skip is used at most once.

Each rejection is paired with a *distinct* earlier skip job, hence the set of matched (skip, rejection) pairs is a collection of *disjoint pairs of jobs*. Therefore, we have  $2R_D \leq n$  (each pair consumes two distinct jobs from the  $n$  arrivals), thereby  $R \leq \lfloor n/2 \rfloor$ . Since  $n \leq m$ , the stated bound  $R_D(M_{CFM}) \leq \lfloor n/2 \rfloor$  follows.

The bound is tight up to the floor: under strictly anti-sorted costs of Definition 7, an adversary can realize  $\lfloor m/2 \rfloor$  units of regret by partitioning providers into disjoint pairs and, for each pair, presenting two jobs in order: an “easy” job feasible for both, then a “threshold” job feasible only for the right provider  $2j$ . **GSM** matches both jobs in each pair, while **CFM** matches exactly one, yielding one unit of regret per pair.  $\square$

Therefore, we conclude that in the case of  $m$  providers, the worst-case D-regret for **CFM** is bounded above by  $\lfloor n/2 \rfloor$ , where  $n$  is the number of online arriving jobs.

To conclude, we have shown that for any job arrival sequence, **CFM** as a incentive-compatible and feasibility-aware matching algorithm, yields a competitive ratio of at least  $1/2$  compared to its optimal counterpart **GSM** and **GCM**.

In Appendix 10.2, we extend the analysis to the multi-period case with two providers, where over an evaluation period  $T_{eval} = \gcd(\tau_s(0), \tau_l(0))$  with  $\tau_s(0) < \tau_l$ , **CFM** induces at most one more infeasible match, hence  $O(1)$ , which is sublinear in the horizon length with a constrained adversary. In the case of multi-period multi-providers, we leave as an open question whether, under aggregate workload constraints, the two-provider  $O(1)$ -per-evaluation-period bound between **CFM** and **GSM** lifts to  $m$  providers, to start, we provide a generalized lemma on a constrained adversary’s optimal strategy against **CFM** and **GSM**.

## 7. Incomplete-Information Extensions

### 7.1. Malicious/Lazy Provider with Verifier

In addition to the complete-information setting where all providers and job submitters are honest, we should also consider cases where the providers are lazy or dishonest. Therefore, for every job that is completed, we initiate a verification game as described in detail by Arun et al. (2025).

### 7.2. Dishonest Provider and Unknown Job Length

The current model can be extended to incorporate an incomplete-information setup when (1) the actual completion time of a job is known to the user, and (2) when the providers are not honest about their capacity (completion time, machine type). To remedy the incomplete market and allow for this mechanism to work, we only need two simple assumptions on the provider side: (1) there are at least  $\frac{1}{n}$  providers who are honest, and (2) job length is bounded by  $w_{\max}$ .



**Environment.** A job has size  $w \sim f_w$  with support  $[0, w_{\max}]$ . Providers  $s_i$  have staking time  $\tau_i$ , registered minimum willing hourly price  $\hat{c}_i$  and market price  $P$ . Feasible providers are  $\mathcal{F} := \{s_i : \tau_{s_i} \geq w_{\max}\}$ . For job  $w$ , provider  $s_i$  would complete in time  $t_i(w) \in \mathbb{R}_+$ . Fix  $n \in \mathbb{N}$  and let  $S_n \subseteq \mathcal{F}$  be the  $n$  cheapest feasible providers by  $\hat{c}_i \leq P$  (ties arbitrary). Each  $s_i \in S_n$  forms a private estimate  $\hat{t}_i$  using the released meta-data (and any private proxy signals); no proxy is disclosed.

**Mechanism (provider's race with stake and tolerance).** Fix a tolerance  $\varepsilon \geq 0$  and a stake amount  $B$  strictly larger than the task reward, i.e.,

$$B > \max_{i \in S_n} P \cdot \hat{t}_i$$

1. *Pre-selection:* Form  $S_n$ .
2. *Bids and stake:* Each  $i \in S_n$  submits a single quote  $\hat{t}_i$  and locks the stake  $B$ .
3. *Allocation and payment:* Let  $i^* \in \arg \min_{i \in S_n} \hat{t}_i$  and award the job to  $s_i^*$  at price  $P^t$ .
4. *Delivery rule (tolerance window):* Completion is deemed on-time if and only if

$$|t_{i^*}(w) - \hat{t}_{i^*}| \leq \varepsilon.$$

If on-time, unlock  $B$  to  $i^*$ ; otherwise (i.e., not finished within the tolerance window), the entire stake  $B$  is forfeited. Finishing early is tolerated and will be paid the full reported length at the market price.

**Proposition 10** (selection optimality and  $\varepsilon$ -accurate quote). *With  $B > \max_{i \in S_n} p_i$  and risk-neutral providers, define  $t^*(w) := \min_{i \in S_n} t_i(w)$  and  $q(w) := \min_{i \in S_n} \hat{t}_i$ , any weakly undominated strategy for provider  $i$  satisfies*

$$\hat{t}_i \in [t_i(w) - \varepsilon, t_i(w) + \varepsilon].$$

Consequently, the mechanism selects a provider whose true completion time is the fastest among the  $n$  racers,

$$t_{i^*}(w) = t^*(w),$$

and the market quote is  $\varepsilon$ -accurate:

$$q(w) \in [t^*(w) - \varepsilon, t^*(w) + \varepsilon].$$

In particular, if all providers best respond by choosing the lowest undominated quote,  $\hat{t}_i = t_i(w) - \varepsilon$ , then  $q(w) = t^*(w) - \varepsilon$ .

*Proof.* If  $s_i$  reports  $\hat{t}_i < t_i(w) - \varepsilon$  and wins, lateness beyond the tolerance is certain and the entire stake  $B > p_i$  is forfeited, yielding a strictly negative payoff; projecting the report up to  $t_i(w) - \varepsilon$  weakly increases expected utility against any opponents' bids. If  $i$  reports  $\hat{t}_i > t_i(w) + \varepsilon$ , lowering the report down to  $t_i(w) + \varepsilon$  weakly raises the win probability without affecting feasibility or payoffs upon winning. Hence reports outside  $[t_i(w) - \varepsilon, t_i(w) + \varepsilon]$  are weakly dominated. Because each  $i$  can (and in equilibrium does) choose the smallest undominated report, the ordering of  $\{\hat{t}_i\}$  coincides with the ordering of  $\{t_i(w)\}$ , so the winner has  $t_{i^*}(w) = \min_{i \in S_n} t_i(w)$ . The bound on  $q(w)$  follows from  $\hat{t}_i \in [t_i(w) - \varepsilon, t_i(w) + \varepsilon]$ .  $\square$   $\square$

## 8. Discussions

### 8.1. Floor Price Updating for Market Healthiness

Recall that each provider registers the minimum price  $\hat{c}_i$  at which they are willing to supply, and the market maintains a floor price  $P_f^t$ . We assume a reduced form pricing function, such that AMM sets the



instantaneous market price as

$$P^t = \begin{cases} P_f^t, & D^t \leq S_f^t, \\ P_f^t \cdot f\left(\frac{D^t}{S_f^t}\right), & D^t > S_f^t, \end{cases}$$

where  $D^t$  is the demand and  $S_f^t$  is the total supply eligible at or below the floor. This ensures that when demand is moderate, jobs transact at the floor, creating stable and predictable pricing, while excess demand triggers price escalation. Thus, the *healthiness* of the floor depends on how well it reflects the actual cost of mobilizing marginal supply. We evaluate this from two perspectives: (i) *price stability*, meaning the floor captures enough eligible supply to prevent frequent spikes; and (ii) *budget feasibility*, meaning most users can lock jobs at the floor.

To update the floor in practice, we use the observed costs of the *marginal providers*, i.e., the most expensive matched providers in each period. Formally, let  $\mathcal{S}^{t*} \subseteq \mathcal{S}^t$  denote the set of matched providers at time  $t$ . Define the period's marginal cost as

$$\max\{\hat{c}_i \mid s_i \in \mathcal{S}^{t*}\}.$$

The updated floor is then set as the time average of these marginal costs over the past  $T$  periods:

$$P_f^t = \frac{1}{T} \sum_{\tau=t-T+1}^t \max\{\hat{c}_i \mid s_i \in \mathcal{S}_\tau^*\}.$$

**Rationale.** The core idea is to anchor the floor price to the *empirical marginal cost of supply*, i.e., the price of the most expensive provider who was actually needed to serve demand. If this marginal cost is consistently *below* the floor, then the floor is set too high: the market is counting too many providers as eligible at floor, demand is always satisfied cheaply, and the system remains over-supplied at the floor price. If the marginal cost is consistently *above* the floor, then the floor is set too low: many periods require bringing in higher-cost providers, so the AMM often lifts the price above the floor, creating frequent and destabilizing spikes.

By updating the floor to the time-averaged marginal cost over the past  $T$  periods, the mechanism automatically corrects either imbalance. When the floor is too high, the update pulls it downward toward the observed marginal costs; when the floor is too low, the update pushes it upward. Averaging over  $T$  periods smooths out noise from temporary shocks, so that the floor reflects longer-run conditions rather than one-off fluctuations. In this way, the updated floor tracks the true cost of mobilizing supply, ensuring that (i) the market usually clears at the floor without spikes-price stability, (ii) users can reliably plan within predictable budgets- budget feasibility.

## 9. Conclusion and Future Directions

We proposed a decentralized two-sided market design that models compute capacity as a time-bound, perishable commodity, enabled by recent advances in computational reproducibility and verifiability (Arun et al., 2025). Our approach (i) prices unit-period capacity to sidestep the combinatorial complexity of traditional auction-based discovery, and (ii) couples premium sharing with matching to obtain incentive compatibility for both sides under our model. On performance, in the single-period setting with  $m$  providers and  $n$  jobs, the algorithm attains  $1/2$  competitive ratio regret: at most  $\lfloor n/2 \rfloor$  relative to the provider-side optimum and at most  $\lfloor m/2 \rfloor$  relative to the user-side optimum. In the multi-period setting with two providers, viewed as a robustness check against an adaptive adversary, the excess infeasible-match differential between **CFM** and **GSM** is at most one per evaluation window  $T_{eval} = \text{lcm}(\tau_s(0), \tau_l(0))$ , yielding a long-run average gap of at most  $1/L$  that vanishes as capacities scale;



moreover, this gap arises only in the anti-sorted gap-1 timing window when initial capacities are coprime and disappears when  $\gcd(\tau_s(0), \tau_l(0)) > 1$ .

Looking ahead, a natural question is whether the two-provider analysis extends cleanly to  $m$  providers. We leave it as an open problem the development of sharp multi-period bounds under aggregate workload constraints. Another promising direction is endogenous floor design: choosing  $P_f^t$  to minimize  $\sum_t (P^t - P_f^t)$  while preserving participation and incentives, potentially via online control or learning. It is also important to study stochastic arrival models (e.g., renewal or Markov-modulated processes) in place of fully adversarial input, to incorporate heterogeneous jobs and preemption with explicit fragmentation costs, and to analyze strategic robustness in the presence of misreports, collusion, and Sybil behavior. Finally, an implementation-focused agenda—including measurements of verification overhead and latency in verifiable-compute backends, trace-driven simulations or testbeds to quantify welfare and utilization, and fairness or service-quality constraints on both sides—would ground the theory in practice and inform mechanism refinements.



## References

- Gagan Aggarwal, Gagan Goel, Chinmay Karande, and Aranyak Mehta. Online vertex-weighted bipartite matching and single-bid budgeted allocations. In *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pages 1253–1264. SIAM, 2011.
- Miguel F Anjos, Russell CH Cheng, and Christine SM Currie. Optimal pricing policies for perishable products. *European Journal of Operational Research*, 166(1):246–254, 2005.
- Arasu Arun, Adam St Arnaud, Alexey Titov, Brian Wilcox, Viktor Kolobaric, Marc Brinkmann, Oguzhan Ersoy, Ben Fielding, and Joseph Bonneau. Verde: Verification via refereed delegation for machine learning programs. *arXiv preprint arXiv:2502.19405*, 2025.
- Itai Ashlagi, Aviv Cohen, Yash Kanoria, et al. Matching queues, flexibility and incentives. *arXiv preprint arXiv:2006.08863*, 2020.
- Yixin Bao, Yanghua Peng, Chuan Wu, and Zongpeng Li. Online job scheduling in distributed machine learning clusters. In *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*, pages 495–503. IEEE, 2018.
- Aaron L. Bodoh-Creed, Jörn Boehnke, and Brent Hickman. How efficient are decentralized auction platforms? *The Review of Economic Studies*, 88(1):91–125, 2021. doi: 10.1093/restud/rdaa017.
- Gerth Stølting Brodal, Loukas Georgiadis, Kristoffer Arnsfelt Hansen, and Irit Katriel. Dynamic matchings in convex bipartite graphs. In *International Symposium on Mathematical Foundations of Computer Science*, pages 406–417. Springer, 2007.
- Yufeng Cao, Anton Kleywegt, and He Wang. Dynamic pricing for two-sided marketplaces with offer expiration. *Management Science*, 2025.
- Data Bridge Market Research. Global cloud computing market size, share and trends analysis report 2024–2032, 2024. URL <https://www.databridgemarketresearch.com/reports/global-cloud-computing-market>. The global cloud computing market size was valued at USD 557.66 B in 2024 and is projected to reach USD 1,705.89 B by 2032.
- Yongheng Deng, Feng Lyu, Ju Ren, Huaqing Wu, Yuezhi Zhou, Yaoxue Zhang, and Xuemin Shen. Auction: Automated and quality-aware client selection framework for efficient federated learning. *IEEE Transactions on Parallel and Distributed Systems*, 33(8):1996–2009, 2021.
- Ludwig Dierks and Sven Seuken. Cloud pricing: The spot market strikes back. *Management Science*, 68(1):105–122, 2022.
- Shahar Dobzinski and Jan Vondrák. The computational complexity of truthfulness in combinatorial auctions. In *Proceedings of the 13th ACM Conference on Electronic Commerce*, pages 405–422, 2012.
- Youyi Feng and Guillermo Gallego. Perishable asset revenue management with markovian time dependent demand intensities. *Management science*, 46(7):941–956, 2000.
- Nicolás Figueroa and Vasiliki Skreta. The role of outside options in auction design. 2007.
- Walter Frick and Felix Salmon. The ai boom needs a market for compute — just like oil and spectrum. *Bloomberg*, September 28 2025. URL <https://www.bloomberg.com/news/newsletters/2025-09-28/artificial-intelligence-needs-a-market-for-compute>. Newsletter: The Forecast, UTC time 08:00 AM.
- Eugene Furman and Adam Diamant. Optimal capacity planning for cloud service providers with periodic, time-varying demand. *Time-Varying Demand (July 10, 2020)*, 2020.
- Guillermo Gallego and Ming Hu. Dynamic pricing of perishable assets under competition. *Management Science*, 60(5):1241–1259, 2014.
- Eyran J Gisches, Hang Qi, William J Becker, and Amnon Rapoport. Strategic retailers and myopic consumers: Competitive pricing of perishable goods. *Journal of Behavioral and Experimental Economics*, 92:101700, 2021.



- Fred Glover. Maximum matching in a convex bipartite graph. *Naval research logistics quarterly*, 14(3):313–316, 1967.
- Shuai Gong, Xiaoyu Yin, Min Zhang, and Wanwan Cao. Improvement of cloud platform utilization based on evaluation and optimization of computing resource runtime surplus capacity. In *International Conference on Frontier Computing*, pages 160–166. Springer, 2023.
- Grand View Research. Cloud computing market size, share & trends analysis report 2024–2030, 2024. URL <https://www.grandviewresearch.com/industry-analysis/cloud-computing-industry>. Global cloud computing market size estimated at USD 752.44 B in 2024 and projected to reach USD 2,390.18 B by 2030.
- Varun Gupta. Greedy algorithm for multiway matching with bounded regret. *Operations Research*, 72(3):1139–1155, 2024.
- Matthew Harding, Kyle Kettler, Carlos Lamarche, and Lala Ma. The (alleged) environmental and social benefits of dynamic pricing. *Journal of Economic Behavior & Organization*, 205:574–593, 2023.
- Darrell Hoy, Nicole Immorlica, and Brendan Lucier. On-demand or spot? selling the cloud to risk-averse customers. In *Web and Internet Economics: 12th International Conference, WINE 2016, Montreal, Canada, December 11–14, 2016, Proceedings 12*, pages 73–86. Springer, 2016.
- Ming Hu and Yun Zhou. Dynamic matching in a two-sided market. Available at SSRN, 2015.
- Zhiyi Huang, Zhihao Gavin Tang, and David Wajc. Online matching: A brief survey. *ACM SIGecom Exchanges*, 22(1):135–158, 2024.
- Paul L Joskow and Catherine D Wolfram. Dynamic pricing of electricity. *American Economic Review*, 102(3):381–385, 2012.
- Bruno Jullien, Alessandro Pavan, and Marc Rysman. Two-sided markets, pricing, and network effects. In *Handbook of industrial organization*, volume 4, pages 485–592. Elsevier, 2021.
- Richard M Karp, Umesh V Vazirani, and Vijay V Vazirani. An optimal algorithm for on-line bipartite matching. In *Proceedings of the twenty-second annual ACM symposium on Theory of computing*, pages 352–358, 1990.
- Ian A Kash, Peter Key, and Warut Suksompong. Simple pricing schemes for the cloud. *ACM Transactions on Economics and Computation (TEAC)*, 7(2):1–27, 2019.
- Vijay Krishna and Motty Perry. Efficient mechanism design, 1998.
- Chien-Yu Lai. *Pricing on Two-Sided Market of Ride-Sharing Platform*. The University of Chicago, 2021.
- SangMok Lee. Incentive compatibility of large centralized matching markets. *Working Paper, University of Pennsylvania*, 2016.
- Kevin Leyton-Brown, Eugene Nudelman, and Yoav Shoham. Learning the empirical hardness of optimization problems: The case of combinatorial auctions. In *International conference on principles and practice of constraint programming*, pages 556–572. Springer, 2002.
- Dilip Mookherjee. Decentralization, hierarchies, and incentives: A mechanism design perspective. *Journal of Economic Literature*, 44(2):367–390, 2006.
- Huijie Peng, Yan Cheng, and Xingyuan Li. Real-time pricing method for spot cloud services with non-stationary excess capacity. *Sustainability*, 15(4):3363, 2023.
- Sebastian Perez-Salazar, Ishai Menache, Mohit Singh, and Alejandro Toriello. Dynamic resource allocation in the cloud with near-optimal efficiency. *arXiv preprint arXiv:1809.02688*, 2018. URL <https://arxiv.org/abs/1809.02688>.
- G Vinu Prasad, Abhinandan S Prasad, and Shrisha Rao. A combinatorial auction mechanism for multiple resource procurement in cloud computing. *IEEE Transactions on Cloud Computing*, 6(4):904–914, 2016.



Jean-Charles Rochet and Jean Tirole. Two-sided markets: a progress report. *The RAND journal of economics*, 37(3):645–667, 2006.

Przemyslaw Rys and Maciej Sobolewski. Two-sided platforms: dynamic pricing and multiple equilibria. Technical report, JRC Digital Economy Working Paper, 2020.

Shifeng Shang, Jinlei Jiang, Yongwei Wu, Zhenchun Huang, Guangwen Yang, and Weimin Zheng. Dabgpm: A double auction bayesian game-based pricing model in cloud market. In *Network and Parallel Computing: IFIP International Conference, NPC 2010, Zhengzhou, China, September 13–15, 2010. Proceedings*, pages 155–164. Springer, 2010.

Shaun Sweeney, Chris King, Mark O’Malley, and Robert Shorten. Embracing fairness in consumer electricity markets using an automatic market maker. *arXiv preprint arXiv:2407.20814*, 2024.

Sushil Mahavir Varma, Pornpawee Bumpensanti, Siva Theja Maguluri, and He Wang. Dynamic pricing and matching for two-sided queues. *Operations Research*, 71(1):83–100, 2023.

Yizhao Wang. Two sided market: Price competition in the food delivery market. 2022.

Dongkuo Wu, Xueyi Wang, Xingwei Wang, Min Huang, and Zhitong Wang. Truthful double auction-based resource allocation mechanisms for latency-sensitive applications in edge clouds. In *Wireless Artificial Intelligent Computing Systems and Applications (WASA 2024)*, Lecture Notes in Computer Science, Vol. 14999, page 39–51. Springer, Cham, 2024.

Yuchen Zhang, Long Luo, Gang Sun, Hongfang Yu, and Bo Li. Deadline-aware online job scheduling for distributed training in heterogeneous clusters. *IEEE Transactions on Cloud Computing*, 2025.



## 10. Appendix

### 10.1. Complexity Analysis

In this section, we provide the detailed complexity analysis for each algorithm including the assumptions of data structure, per job time complexity, and space complexity. Also, we present the pseudocodes for **GCM**, **GSM**, and **CFM**.

Algorithm	What it chooses	Data structure (online)	Per-job time	Space
<b>GCM</b> (Greedy Cheapest)	Provider with the smallest cost $\hat{c}_i$ (ignores $\tau_i$ ).	Bucket queue over integer costs $c_i = 10\hat{c}_i \in \{0, \dots, C_{\max}\}$ , where $C_{\max} = \lfloor 10\hat{c}_{\max} \rfloor$ is induced by the current market price. Maintain an array of queues $B[0, \dots, C_{\max}]$ and a pointer to the lowest nonempty bucket.	Pop cheapest: $O(1)$ amortized. Insert: $O(1)$ amortized. Occasional extension when $\hat{c}_{\max}$ grows is amortized $O(1)$ .	$O(n + C_{\max})$
<b>GSM</b> (Greedy Shortest Feasible)	Smallest $\tau_i$ such that $\tau_i > w_j$ ; if none exists, fall back to the provider with the largest $\tau_i$ .	Balanced BST (ordered multiset) keyed by $\tau_i$ . Each key stores a small queue of provider IDs. Maintain the current maximum key $\tau_{\max}$ for the fallback.	Upper_bound search in $\tau$ : $O(\log n)$ . Delete / erase key: $O(\log n)$ . Fallback to $\tau_{\max}$ is $O(1)$ if tracked (otherwise $O(\log n)$ ).	$O(n)$
<b>CFM</b> (Cheapest Feasible Matching)	Among providers with $\tau_i > w_j$ , pick the smallest cost $\hat{c}_i$ ; if none exists, fall back to the provider with the largest $\tau_i$ .	Segment tree over compressed $\tau$ values $\tau^{(1)} < \dots < \tau^{(U)}$ . Each leaf $u$ holds a min-heap $H_u$ of $(c, i)$ for providers at $\tau^{(u)}$ . Every tree node stores the best triple $(c, i, u)$ in its subtree and a boolean “occupied” flag, which also allows a search for the rightmost nonempty leaf (largest $\tau$ ) for fallback.	Range-min query on $\{\tau_i > w_j\}$ : $O(\log U)$ . Delete from a leaf heap: $O(\log k_{\tau})$ . Update path in the tree: $O(\log U)$ . Fallback (rightmost nonempty leaf): $O(\log U)$ . Overall worst-case per job: $O(\log n)$ .	$O(n + U)$

**Table 1** Complexity summary for the three online matching rules. Costs are discrete:  $\hat{c}_i \in \{0, 0.1, 0.2, \dots\}$  so  $c_i = 10\hat{c}_i$  is an integer. GCM uses the current market price  $\hat{c}_{\max}$  as a cap  $C_{\max}$  for the bucket queue. GSM and CFM fall back to the provider with the largest  $\tau_i$  when no provider satisfies  $\tau_i > w_j$ .



---

**Algorithm 1** GCM: Greedy Cheapest Matching

---

```

1: State: Buckets  $B[0, \dots, C_{\max}]$  (queues of provider IDs); pointer  $p$  to lowest nonempty index.
2: Build: For each provider  $(i, \hat{c}_i)$ :
3:    $c_i \leftarrow 10\hat{c}_i$ ; ensure  $c_i \leq C_{\max}$  (extend array if needed); push  $i$  into  $B[c_i]$ ;  $p \leftarrow \min(p, c_i)$ .
4: InsertProvider $(i, \hat{c}_i)$ :
5:    $c_i \leftarrow 10\hat{c}_i$ ; ensure  $c_i \leq C_{\max}$ ; push  $i$  into  $B[c_i]$ ;  $p \leftarrow \min(p, c_i)$ .
6: Match $(w_j)$ :  $\{w_j \text{ unused}\}$ 
7:   while  $p \leq C_{\max}$  and  $B[p]$  empty do  $p \leftarrow p + 1$  end while
8:   if  $p > C_{\max}$  then return FAIL
9:    $i^* \leftarrow \text{pop\_front}(B[p])$ 
10:  return  $i^*$ 

```

---



---

**Algorithm 2** GSM: Greedy Shortest Feasible

---

```

1: State: Balanced BST  $T$  keyed by  $\tau$ ; each key stores a queue of provider IDs.
2: Maintain  $\tau_{\max}$  as the current largest key in  $T$ .
3: Build: For each provider  $(i, \tau_i)$ :
4:   push  $i$  into  $T[\tau_i]$ ; update  $\tau_{\max}$ .
5: InsertProvider $(i, \tau_i)$ :
6:   push  $i$  into  $T[\tau_i]$ ; update  $\tau_{\max}$ .
7: Match $(w_j)$ :
8:    $\text{it} \leftarrow T.\text{UPPER\_BOUND}(w_j)$  {first key with  $\tau > w_j$ }
9:   if  $\text{it} \neq T.\text{END}()$  then
10:     $\tau^* \leftarrow \text{KEY}(\text{it})$ 
11:   else if  $T$  is empty then
12:     return FAIL
13:   else
14:      $\tau^* \leftarrow \tau_{\max}$  {fallback to longest  $\tau$ }
15:   end if
16:    $i^* \leftarrow \text{pop\_front}$  from queue at  $T[\tau^*]$ 
17:   if that queue becomes empty then
18:     erase key  $\tau^*$  from  $T$ ; update  $\tau_{\max}$ 
19:   end if
20:   return  $i^*$ 

```

---




---

**Algorithm 3** CFM: Cheapest Feasible Matching
 

---

```

1: Preprocessing: Let  $\tau^{(1)} < \dots < \tau^{(U)}$  be the sorted distinct staking times. Define index map
    $u = \text{IDX}(\tau)$  so  $\tau^{(u)}$  is the  $u$ -th distinct value.
2: For each  $u$ , maintain a min-heap  $H_u$  of pairs  $(c, i)$  for providers with  $\tau_i = \tau^{(u)}$ .
3: Build a segment tree over indices  $1, \dots, U$ . Each node stores:
4:   best =  $(c, i, u)$ , the minimum-cost provider in its subtree (use  $c = +\infty$  if empty),
5:   occ  $\in \{\text{true}, \text{false}\}$ , indicating whether the subtree is nonempty.
6: LeafUpdate( $u$ ):
7:   if  $H_u$  nonempty then
8:     let top of  $H_u$  be  $(c, i)$ ; set leaf.best =  $(c, i, u)$ , leaf.occ = true
9:   else set leaf.best.c =  $+\infty$ , leaf.occ = false
10:  Propagate updates from leaf to root.
11: Build: For each provider  $(i, \tau_i, \hat{c}_i)$ :
12:    $u \leftarrow \text{IDX}(\tau_i)$ ;  $c \leftarrow 10\hat{c}_i$ ; push  $(c, i)$  into  $H_u$ ; LeafUpdate( $u$ ).
13: InsertProvider( $i, \tau_i, \hat{c}_i$ ):
14:    $u \leftarrow \text{IDX}(\tau_i)$ ;  $c \leftarrow 10\hat{c}_i$ ; push  $(c, i)$  into  $H_u$ ; LeafUpdate( $u$ ).
15: RightmostNonempty():
16:   Starting at root, at each node prefer the right child if its occ is true; otherwise go left. If no such
      child exists, return NONE. This takes  $O(\log U)$  time.
17: Match( $w_j$ ):
18:   Find  $t$  such that  $\tau^{(t)} \leq w_j < \tau^{(t+1)}$  by binary search (if all  $\tau^{(u)} > w_j$ , set  $t = 0$ ; if all  $\tau^{(u)} \leq w_j$ , set
       $t = U$ ).
19:   Query segment tree on interval  $(t + 1, \dots, U)$  to obtain  $q.\text{best} = (c, i, u)$ .
20:   if  $c < +\infty$  then {feasible cheapest exists}
21:      $(\_, i^*, u^*) \leftarrow (c, i, u)$ 
22:   else {fallback to provider with largest  $\tau$ }
23:      $u^* \leftarrow \text{RightmostNonempty}()$ 
24:     if  $u^* = \text{NONE}$  then return FAIL
25:      $i^* \leftarrow \text{top}(H_{u^*}).i$ 
26:   end if
27:   Pop the top of  $H_{u^*}$ ; LeafUpdate( $u^*$ )
28:   return  $i^*$ 

```

---

## 10.2. Multi-period Robust Analysis

Having established tight single-period bounds on the algorithmic differential between **GSM** and **CFM**, we now turn to dynamics. In the multi-period setting, providers' remaining staking times evolve, so today's assignment shapes tomorrow's feasibility. Rather than posit a multi-period oracle and define regret to it, we treat the multi-period results as a robustness check: we compare **GSM** (**GCM**) and **CFM** under an adaptive adversary and measure their dynamic performance differential in accumulated infeasibilities (cost) over an evaluation horizon. To make the mechanisms transparent, we focus first on the minimal non-trivial case of two providers and max two jobs per period. This reduction is not ad hoc: it isolates the canonical timing externality where the “anti-sorted” window where the algorithms diverge, allows exact tracking over the natural hyper-period which is the least common multiple of staking times, and produces closed-form conditions that later inform the multi-provider discussion. We leave as an open question whether these pairwise insights admit a sharp decomposition for arbitrary numbers of providers under aggregate workload constraints, e.g. additive or near-additive bounds across interacting pairs, and whether such bounds are tight.

For multi-period case with 2 providers with initial staking time  $\tau_l(0), \tau_s(0)$ , there exists a deterministic



time  $T$  such that all providers restake, therefore we will evaluate an adaptive adversary's infeasible match maximization to **CFM** and **GSM** over the period of  $T_{eval}(\tau_s(0), \tau_l(0)) = \text{lcm}(\tau_s(0), \tau_l(0))$ .

Before beginning the robustness analysis, we give the general formulation of this multi-period bipartite matching problem with an adaptive adversary in addition to the single period rules described in the previous section.

- Players' staking time  $\tau_s(t)$  evolves as

$$\tau_s(t+1) = \begin{cases} \tau_s(t) - 1 & \text{if } \tau_s(t) \geq 2 \\ \tau_s(0) & \text{else} \end{cases}$$

- Unfinished jobs must be released in the period when the infeasible provider finishes, but the adversary can change the order.
- The adaptive adversary is constrained in a way that
  - The adversary has to send the same number of jobs to the number of idle providers.
  - For period  $t$ , denote  $w_j^t$  as the  $j$ -th scheduled job in that period, we have  $\sum_j w_j^t \leq \sum_{s \text{ is idle}} \tau_s(t)$ .
  - Any infeasible match that results in a residual job, that job will have to be released at the instant when next provider(s) become available and count as one of the released job(s) in that period.

The reason we have a constrained adversary is that, in a finite period setting, there always exists a universal adversary strategy with job length equal to  $\text{lcm}(\tau_s(0), \tau_l(0)) + 1$ , which maximize the number of infeasible matches irrespective matching algorithm used; because at each restake, the idle provider will be paired with an infeasible residual until the end of  $T_{eval}(\tau_s(0), \tau_l(0))$ . Therefore, to assess the robustness of **CFM** and **GSM** with the two providers' type distribution being anti-sorted (Definition 7)

**Lemma 5.** *A constrained adversary can create at most one infeasible match per period.*

*Proof.* Since at any time  $t$  with two providers,  $w_{t,1} + w_{t,2} \leq \tau_s(t) + \tau_l(t)$  and  $w_{t,i} \leq \max(\tau_s(t), \tau_l(t))$ , no infeasible match can be enforced on the provider with longer remaining time. Without loss of generality, assume  $w_{t,1}$  enforces an infeasible match at  $t$  on the short provider, such that  $w_{t,1} > \min(\tau_l(t), \tau_s(t))$ , then  $\max(w_{t,2}) < \max(\tau_l(t), \tau_s(t))$ , which cannot enforce an infeasible match on the remaining long provider.  $\square$

We now give lemmas to pin down the sufficient conditions where each algorithm can create an infeasible match on one of the providers.

**Lemma 6.** *At most one unit of infeasible match can be enforced during a single restake cycle of any provider.*

*Proof.* One infeasible match requires assigning a job longer than the provider's *current* remaining time at an idle instant. If at some idle time within a restake cycle the remaining time is already strictly less than the full stake  $\tau_i(0)$ , then an infeasible match must have been enforced earlier in that cycle (or no infeasible match is possible), so at most one such infeasible match can occur per cycle.  $\square$

**Lemma 7** (Sufficient condition for an infeasible match on **CFM**). *A constrained adversary can create an infeasible match against **CFM** if:*

- *The providers are anti-sorted, such that  $\tau_l(t) > \tau_s(t)$ , then it is always possible to create an infeasible match.*



- The providers are sorted, such that  $\tau_l(t) \leq \tau_s(t)$  then it is only possible to create an infeasible match when  $\tau_s(t) - \tau_l(t) \geq 2$ .

*Proof.* If the providers are anti-sorted, then the optimal adversary strategy is to send a 1-period job first, which will be matched with the currently long provider, and then send a job  $\tau_s(t) \leq w_j^t \leq \tau_l(t)$ , on the short provider to create an infeasible match. If the providers are sorted, then to keep  $\tau_s(t)$  busy, first send  $w_{t,1} = \tau_l(t) + 1$  to keep it busy, and then stack the rest to the long provider to enforce the regret. Therefore, we must have  $\tau_s(t) - \tau_l(t) \geq 2$  when the providers are sorted.  $\square$

**Lemma 8** (Sufficient Condition for infeasible match on **GSM**). *An adversary can create an infeasible match against **GSM** if  $\max(\tau_s(t), \tau_l(t)) - \min(\tau_s(t), \tau_l(t)) \geq 2$ .*

*Proof.* Since when the providers are sorted, **CFM** is equivalent to **GSM** according to Lemma 4, therefore the sufficient condition for the adversary to create an infeasible match on **GSM** is the same as **CFM** when the providers are sorted.  $\square$

From the two lemmas above, we can conclude the following lemma on when the sufficient condition for an infeasible match is different between **CFM** and **GSM**.

**Lemma 9.** *If  $\tau_l(t) - \tau_s(t) = 1$  and  $\tau_l(t) > \tau_s(t) \geq 2$ , then the adversary can enforce an infeasible match on **CFM** but not on **GSM**.*

*Proof.* Under  $\tau_l(t) - \tau_s(t) = 1$ , if the **GSM** adversary want to enforce a infeasible match on the short provider, she has to send  $w_{t,1} = \tau_s(t) + 1$  to occupy the long provider first according to Lemma 8, by doing so, her second job's length is constrained by  $w_{t,2} \leq \tau_s(t) + \tau_l(t) - (t) - 1 = \tau_s(t)$ , which is not enough to enforce a regret.

When going against **CFM**, the adversary can send  $w_{t,1} = 1$  to occupy the long provider, and then send  $w_{t,2} = \tau_l(t) + \tau_s(t) - 1 > \tau_s(t)$  to the shorter provider.  $\square$

With these results in hand, we next characterize when the anti-sorted gap 1 window occurs along the deterministic  $T_{eval}(\tau_s(0), \tau_l(0))$ .

**Proposition 11.** *Let  $\tau_l(0) > \tau_s(0)$  be positive integers and define the sequences*

$$\tau_l(t) \equiv -t \pmod{\tau_l(0)}, \quad \tau_s(t) \equiv -t \pmod{\tau_s(0)},$$

*with values taken in  $\{1, \dots, \tau_l(0)\}$  and  $\{1, \dots, \tau_s(0)\}$ , respectively.*

1. *If  $\gcd(\tau_l(0), \tau_s(0)) > 1$ , there is no  $t$  such that  $\tau_l(t) - \tau_s(t) = 1$ .*
2. *If  $\gcd(\tau_l(0), \tau_s(0)) = 1$ , then in each period of length  $L$  there is exactly one contiguous block of length  $\tau_s(0)$  of times  $t$  such that  $\tau_l(t) - \tau_s(t) = 1$ .*

*Proof.* Let  $\tau_l(0) > \tau_s(0)$  be positive integers. Define the sequences

$$\tau_l(t), \tau_s(t) \in \{1, 2, \dots, \tau_l(0)\} \times \{1, 2, \dots, \tau_s(0)\}, \quad t \in \mathbb{Z},$$

for each step,

$$\tau_l(t+1) = \begin{cases} \tau_l(t) - 1, & \tau_l(t) > 1, \\ \tau_l(0), & \tau_l(t) = 1, \end{cases} \quad \tau_s(t+1) = \begin{cases} \tau_s(t) - 1, & \tau_s(t) > 1, \\ \tau_s(0), & \tau_s(t) = 1. \end{cases}$$



Equivalently,

$$\tau_l(t) \equiv -t \pmod{\tau_l(0)}, \quad \tau_s(t) \equiv -t \pmod{\tau_s(0)},$$

with residues interpreted in  $\{1, \dots, \tau\}$  (so that residue 0 corresponds to  $\tau$ ).

Then, the pair  $(\tau_l(t), \tau_s(t))$  is  $\text{lcm}(\tau_l(0), \tau_s(0))$ -periodic.

**Existence of times with  $\tau_l(t) - \tau_s(t) = 1$ .** Fix  $a \in \{1, \dots, \tau_l(0)\}$  and  $b \in \{1, \dots, \tau_s(0)\}$ . The condition

$$\tau_l(t) = a, \quad \tau_s(t) = b$$

is equivalent to the simultaneous congruences

$$t \equiv -a \pmod{\tau_l(0)}, \quad t \equiv -b \pmod{\tau_s(0)}. \tag{*}$$

By the Chinese remainder theorem,  $(*)$  has a solution if and only if

$$-a \equiv -b \pmod{g} \iff a \equiv b \pmod{g},$$

where  $g = \gcd(\tau_l(0), \tau_s(0))$ .

Now suppose  $a - b = 1$ .

- If  $g > 1$ , then  $a \equiv b \pmod{g}$  cannot hold, because  $1 \not\equiv 0 \pmod{g}$ . Hence there are no solutions: no time  $t$  satisfies  $\tau_l(t) - \tau_s(t) = 1$ .
- If  $g = 1$ , then the congruences always have a unique solution modulo  $\text{lcm}(\tau_l(0), \tau_s(0))$ . Thus such  $t$  do exist.

**Structure of the set of solutions when  $\gcd(\tau_l(0), \tau_s(0)) = 1$ .** Consider the set

$$\mathbf{S} := \{(a, b) : 1 \leq b \leq \tau_s(0), a = b + 1\}.$$

Since  $\tau_l(0) > \tau_s(0)$ , all such pairs are valid. Thus  $|\mathbf{S}| = \tau_s(0)$ .

Because  $\gcd(\tau_l(0), \tau_s(0)) = 1$ , the map  $t \mapsto (\tau_l(t), \tau_s(t))$  is a bijection from  $\mathbb{Z}/\text{lcm}(\tau_l(0), \tau_s(0))\mathbb{Z}$  to the full grid  $\{1, \dots, \tau_l(0)\} \times \{1, \dots, \tau_s(0)\}$ . Therefore, each pair in  $\mathbf{S}$  occurs exactly once per hyper-period. Consequently, there are exactly  $\tau_s(0)$  distinct times with  $\tau_l(t) - \tau_s(t) = 1$  in each period of length  $\text{lcm}(\tau_l(0), \tau_s(0))$ .

**Consecutivity.** Within  $S$ , the transition dynamics are

$$(b+1, b) \mapsto (b, b-1).$$

For  $b \geq 2$ , the image remains in  $S$ ; only  $(2, 1)$  exits  $S$ . Thus the occurrences in  $S$  appear consecutively as

$$(\tau_s(0)+1, \tau_s(0)), (\tau_s(0), \tau_s(0)-1), \dots, (3, 2), (2, 1).$$

This shows that all  $\tau_s(0)$  occurrences appear as a single contiguous block.

□

**Proposition 12.** For two providers with initial staking time  $\tau_l(0), \tau_s(0)$  such that  $\tau_l(0) > \tau_s(0) \geq 2$ , define restaking hyper-period, as  $\text{lcm}(\tau_s(0), \tau_l(0))$ , an adversary can create at most one extra infeasible match against **CFM** than against **GSM** during  $\text{lcm}(\tau_s(0), \tau_l(0))$  if  $\tau_s(0), \tau_l(0)$  are coprime. More specifically, if  $\gcd(\tau_s(0), \tau_l(0)) \geq 2$ , the adversary cannot create more infeasible matches against **CFM**. Thus, the time average infeasible match difference between **CFM** and **GSM** is bounded above by  $\frac{1}{\tau_s(0)\tau_l(0)}$ .



*Proof.* From Lemma 7 and Lemma 8, the only time when a **CFM**-adversary can create an infeasible match but **GSM**-adversary cannot is when providers are anti-sorted, and  $\tau_l(t) - \tau_s(t) = 1$ , from Proposition 11, in one hyper-period  $\text{lcm}(\tau_s(0), \tau_l(0))$ , there are exactly  $\tau_s(0)$  distinct times  $t$  such that  $\tau_l(t) - \tau_s(t) = 1$ , and they occur consecutively. From, Lemma 6, only 1 infeasible match can be created during this consecutive  $\tau_s(0)$ -period. Therefore,  $\text{lcm}(\tau_s(0), \tau_l(0))$  hyper-period, the adversary for **CFM** can create at most **one more** regret compared to **GSM**.

Moreover, if  $\gcd(\tau_s(0), \tau_l(0)) \geq 2$ , according to Proposition 11, no such  $t$  exists such that  $\tau_l(t) - \tau_s(t) = 1$ , therefore, an adversary cannot create more infeasible matches on **CFM** over **GSM**.

Additionally, the time-average infeasible match difference between **CFM** and **GSM** in a two-provider setting is bounded above by  $\frac{1}{\text{lcm}(\tau_s(0), \tau_l(0))}$ .  $\square$

We have now shown that in the case of two providers, over each evaluation hyper-period  $T_{eval}$ , the **CFM-GSM** excess infeasibility is  $O(1)$ , hence sublinear in the horizon length; over time its long-run per-period rate is at most  $1/L$  and vanishes as capacities scale. We leave as an open question whether, under aggregate workload constraints, the two-provider  $O(1)$  -per-hyper-period bound lifts to  $m$  providers case. The above proposition can also be extended to a multi-provider setting, suppose that there are  $m$  providers, since the infeasible match between **CFM** and **GSM** in a two-provider setting is only determined by whether  $\tau_s(0)$  and  $\tau_l(0)$  are coprime, we use induction to show that the same result hold for **GSM** and **CFM** comparison for  $m > 2$  provider case.

**Lemma 10** (Peel-and-load characterization of maximum infeasible matches). *Let there be  $m$  providers with anti-sorted staking times (capacities)*

$$\tau_1 \geq \tau_2 \geq \dots \geq \tau_m \geq 2.$$

*An adversary has total budget  $B = \sum_{i=1}^m \tau_i$  and must assign exactly one job to each provider. A job  $w_i \leq \tau_i$  is feasible, otherwise it causes an infeasible match (an overload). On overload, the matching algorithm always removes the provider with the largest remaining capacity.*

*The maximum number of infeasible matches that an adversary can enforce against either policy follows the peel-and-load inductive process: the adversary first peels the longer providers by sending feasible jobs and saving budget, and then uses the saved budget to load shorter providers by sending infeasible jobs.*

Let  $\text{Save}_i^{(M)}$  denote the total savings after peeling the first  $i$  longest providers under policy  $P$ . Then:

$$\begin{aligned} \text{Save}_i^{(\text{CFM})} &= \sum_{j=1}^i (\tau_j - 1), \\ \text{Save}_i^{(\text{GSM})} &= \sum_{j=1}^i (\tau_j - \tau_{j+1} - 1), \quad \text{where } \tau_{m+1} := 1. \end{aligned}$$

*The minimum number of long providers that must be peeled before the savings are sufficient to overload all remaining providers is given by*

$$i_{\text{CFM}}^* = \min \left\{ i \in [1, m] : \sum_{j=1}^i (\tau_j - 1) \geq m - i \right\}, \quad (11)$$

$$i_{\text{GSM}}^* = \min \left\{ i \in [1, m] : \sum_{j=1}^i (\tau_j - \tau_{j+1} - 1) \geq m - i \right\}. \quad (12)$$



The maximum achievable number of infeasible matches under each policy is then

$$U_{\max}^{(\mathcal{M})} = m - i_p^*, \quad \mathcal{M} \in \{\text{CFM}, \text{GSM}\}.$$

Moreover, a constrained adversary can enforce at most  $m - 1$  infeasible matches against either policy.

*Proof.* We proceed by induction on the number of providers  $m$ . For the base case  $m = 1$ , no feasible job can create overloads elsewhere. Hence  $U_{\max}^{(\mathcal{M})} = 0$ , and the lemma holds trivially.

Assume the lemma holds for all systems with fewer than  $m$  providers. We now show it holds for  $m$  providers.

The adversary operates according to the *peel-and-load* rule:

- (1) **Peeling stage.** The adversary sequentially removes the  $i$  longest providers using feasible jobs  $w_j \leq \tau_j$ . Each peel saves budget

$$\text{save}_j^{(\mathcal{M})} = \tau_j - w_j,$$

where  $w_j$  is chosen minimally to ensure that provider  $j$  is removed under policy  $P$ :

$$\begin{array}{ll} w_j^{(\text{CFM})} = 1, & \text{save}_j^{(\text{CFM})} = \tau_j - 1, \\ w_j^{(\text{GSM})} = \tau_{j+1} + 1, & \text{save}_j^{(\text{GSM})} = \tau_j - \tau_{j+1} - 1. \end{array}$$

The total savings after peeling  $i$  providers is  $\text{Save}_i^{(\mathcal{M})} = \sum_{j=1}^i \text{save}_j^{(\mathcal{M})}$ .

- (2) **Loading stage.** The adversary uses the saved budget to overload the remaining  $m - i$  providers, each requiring 1 additional unit of length beyond their capacity. A sufficient condition to overload all remaining providers is thus

$$\text{Save}_i^{(\mathcal{M})} \geq m - i.$$

The smallest  $i$  satisfying this condition is  $i_p^*$ , as defined in equation 11–equation 12.

We show both sufficiency and necessity for optimality.

**Sufficiency:** If  $\text{Save}_i^{(\mathcal{M})} \geq m - i$ , the adversary can spend one unit of the saved budget to overload each of the  $m - i$  remaining providers, achieving  $m - i$  infeasible matches.

**Necessity:** Suppose instead that the adversary achieves  $m - (i_p^* - 1)$  infeasible matches by peeling fewer than  $i_p^*$  providers, i.e. using savings  $\text{Save}_{i'}^{(\mathcal{M})}$  with  $i' < i_p^*$ . To overload all remaining  $m - i'$  providers, one must have  $\text{Save}_{i'}^{(\mathcal{M})} \geq m - i'$ . However, by the definition of  $i_p^*$ ,

$$\text{Save}_{i'}^{(\mathcal{M})} < m - i',$$

which is a contradiction. Therefore,  $i_p^*$  is the minimal number of peels necessary, and  $U_{\max}^{(\mathcal{M})} = m - i_p^*$  is the maximal achievable infeasible match gap.

**Upper bound.** At least one provider must be peeled to generate any savings, hence  $i_p^* \geq 1$  and  $U_{\max}^{(\mathcal{M})} = m - i_p^* \leq m - 1$ . When  $i_p^* = 1$ , the adversary can overload all remaining providers after peeling only the longest one. This requires

$$S_1^{(\text{CFM})} = \tau_1 - 1 \geq m - 1, \quad \text{and} \quad S_1^{(\text{GSM})} = \tau_1 - \tau_2 - 1 \geq m - 1,$$

that is,  $\tau_1 \geq m$  for **CFM** and  $\tau_1 - \tau_2 \geq m$  for **GSM**. □