

---

# Wait-Less Offline Tuning and Re-solving for Online Decision Making

---

Jingruo Sun<sup>1</sup> Wenzhi Gao<sup>2</sup> Ellen Vitercik<sup>1,3</sup> Yinyu Ye<sup>1,4,5</sup>

## Abstract

Online linear programming (OLP) has found broad applications in revenue management and resource allocation. State-of-the-art OLP algorithms achieve low regret by repeatedly solving linear programming (LP) subproblems that incorporate updated resource information. However, LP-based methods are computationally expensive and often inefficient for large-scale applications. By contrast, recent first-order OLP algorithms are more computationally efficient but typically suffer from weaker regret guarantees. To address these shortcomings, we propose a new algorithm that combines the strengths of LP-based and first-order OLP algorithms. Our algorithm re-solves the LP subproblems periodically at a predefined frequency  $f$  and uses the latest dual prices to guide online decision-making. In parallel, a first-order method runs during each interval between LP re-solves and smooths resource consumption. Our algorithm achieves  $\mathcal{O}(\log(T/f) + \sqrt{T})$  regret and delivers a “wait-less” online decision-making process that balances computational efficiency and regret guarantees. Extensive experiments demonstrate at least 10-fold improvements in regret over first-order methods and 100-fold improvements in runtime over LP-based methods.

## 1. Introduction

Sequential decision-making has garnered significant attention for its utility in guiding optimal strategies in dynamic environments. The goal is to identify effective decisions and policies in environments where knowledge of the sys-

tem continuously accumulates and evolves. Online Linear Programming (OLP) (Agrawal et al., 2014) offers a powerful framework that encapsulates the core principles of sequential decision-making and has been extensively applied to different domains, including resource allocation (Bal-seiro et al., 2022a), online advertising (Mehta et al., 2007), and inventory management (Talluri et al., 2004).

We study an OLP problem where customers arrive sequentially, each requesting a combination of resources and offering a bidding price. The objective is to determine which resource requests to fulfill to maximize revenue while respecting resource constraints. The challenge is that decisions must be made immediately and irrevocably, relying solely on historical data without knowledge of future arrivals. The goal is to minimize regret with respect to the optimal hindsight linear programming (LP) solution.

To guide real-time decision-making, state-of-the-art OLP algorithms estimate optimal *dual prices* and use them to make decisions. These algorithms fall into two main categories: LP-based and first-order methods. Specifically, LP-based methods update dual prices by repeatedly solving linear programs at each time step with all available information so far. However, the substantial computational demands limit their application in time-sensitive settings. First-order methods offer quick, incremental updates to dual prices using gradient information, but generally fall short of achieving the strong regret guarantee of LP-based methods. These trade-offs motivate an open question at the intersection of online learning and decision-making:

*Can we simultaneously achieve*

*low regret and computational efficiency?*

**Our contributions.** We answer this question in the affirmative. We summarize our contributions as follows:

- *Parallel Multi-Stage Framework:* We separate online learning and decision-making into distinct processes that interact via a feedback loop. We re-solve the LP subproblems periodically at a fixed frequency and feed the updated dual price to guide decision-making until the next time that the LP is re-solved. To further enhance efficiency, we apply a first-order method during the initial and final stages, restarting with the most recent learn-

---

<sup>1</sup>Department of Management Science & Engineering, Stanford University <sup>2</sup>Institute for Computational and Mathematical Engineering, Stanford University <sup>3</sup>Department of Computer Science, Stanford University <sup>4</sup>Chinese University of Hong Kong (Shen Zhen) <sup>5</sup>Hong Kong University of Science and Technology. Correspondence to: Jingruo Sun <jingruo@stanford.edu>.

*Proceedings of the 42<sup>nd</sup> International Conference on Machine Learning*, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

ing outcomes. By integrating the LP-based and first-order techniques, our approach leverages their respective strengths and achieves a balance between regret performance and computational costs.

- *Regret Analysis:* We unify the analysis of LP-based and first-order methods by deriving a new performance metric to account for their inter-dependency. Our analysis yields a “spectrum theorem” that bounds regret for any feasible choice of the re-solving frequency:

**Theorem 1.1** (Informal version of Theorem 3.2). *If we re-solve the LPs every  $f$  time steps within horizon  $T$  and apply a first-order method in the initial and final  $f$  steps, we achieve a worst-case regret of  $\mathcal{O}(\log(T/f) + \sqrt{f})$ .*

In particular,  $f = 1$  yields a pure LP-based method with  $\mathcal{O}(\log T)$  regret, while  $f = T$  reduces to a pure first-order method with  $\mathcal{O}(\sqrt{T})$  regret. By choosing an “optimal”  $f$ , one can achieve the best possible regret based on available computational resources and enable a “wait-less” decision-making system across all time steps.

- *Experiments:* Through experiments across diverse distributions, we demonstrate that our algorithm achieves over a 10-fold improvement in regret compared to first-order methods and over a 100-fold improvement in run-time with comparable regret to LP-based methods.

**Key Challenges.** OLP methods face a fundamental challenge in balancing efficient real-time decision-making and accurate dual-price learning. Specifically, LP-based methods ensure high-quality decisions with  $\mathcal{O}(\log T)$  regret, but their computational costs increase quadratically with problem size. A natural alternative—batching customers and solving the LP once every  $f$  arrivals—reduces this cost, but leaves customers waiting until the batch concludes. In contrast, first-order methods update dual prices via gradient information, allowing for faster computation. However, these methods require small step sizes to maintain decision quality, which slows their adaptation to new data. Even using different step sizes for learning and decision-making only achieves  $\mathcal{O}(T^{1/3})$  regret when the distribution of customers has continuous support. Our parallel framework addresses these issues: an LP-based method periodically refines dual prices, while a first-order method immediately processes arriving customers using the most recent dual updates, thereby eliminating delays.

A second challenge emerges when analyzing the regret of this hybrid approach. LP-based methods rely on a stopping-time analysis (halting when resources are depleted), whereas first-order methods track constraint violations as part of the regret. Since these two formulations are fundamentally different, existing regret bounds cannot be

naively combined. We overcome this challenge by introducing a unified performance metric that decomposes regret into three components: dual convergence, constraint violation, and leftover resources. This approach yields the first integrated analysis of LP-based and first-order OLP, culminating in our  $\mathcal{O}(\log(T/f) + \sqrt{f})$  regret bound.

## 1.1. Related Literature

Online resource allocation and OLP problems have been widely studied under two predominant models: the stochastic input (Goel & Mehta, 2008; Devanur et al., 2019) and stochastic permutation models (Agrawal et al., 2014; Gupta & Molinaro, 2016). We study the former, where each customer’s resource request and bidding price are drawn i.i.d. from an unknown distribution. We use expected regret and constraint violation as the performance metric. We summarize some recent works in Table 1.

**LP-based Methods.** These methods solve the OLP dual repeatedly to update dual prices and make decisions. Early approaches enforced a fixed average resource constraint (Agrawal et al., 2014), whereas recent work dynamically tracks remaining resources and enables LP-based methods to achieve  $\mathcal{O}(\log T)$  regret under continuous support and non-degeneracy assumptions (Li & Ye, 2022). Variants include multi-secretary problems (Bray, 2022), regularized resource constraints (Ma et al., 2024), and finite-support distributions yielding constant regret (Chen et al., 2024).

**First-order Methods.** These methods generate decisions using gradient updates without solving LPs, enabling efficient computation. They achieve  $\mathcal{O}(\sqrt{T})$  regret with mirror descent (Li et al., 2020; Balseiro et al., 2022a) and  $\mathcal{O}(T^{3/8})$  under finite support and non-degeneracy assumptions (Sun et al., 2020). Variants include proximal updates (Gao et al., 2023), momentum-based mirror descent (Balseiro et al., 2022b), resource adjustments (Ma et al., 2024), and restart strategies yielding  $\mathcal{O}(T^{1/3})$  regret (Gao et al., 2024).

**Delay in Decision-making.** Delays arise from the time-consuming process of solving the large-scale, up-to-date LP subproblems for each customer. Golrezaei & Yao (2021) study a mix of impatient and partially patient customers, while Xie et al. (2023) show that batching requests can reduce regret. Concurrent work (Xu et al., 2024) reduces delay by solving LPs in batches but assumes a lower bound on resource requests and still requires waiting in initial and final batches. To achieve delay-free decisions, we shift the re-solving process offline and only fine-tune the solution online. Compared with previous works, our framework imposes standard assumptions, achieves  $\mathcal{O}(\log(T/f) + \sqrt{f})$  regret, and ensures “wait-less” decision-making throughout the entire horizon.

Table 1: Performances of Dual Algorithms in Recent OLP Literature

Paper	Setting	Algorithm	Regret	Decision-Making
(Li & Ye, 2022)	Bounded, continuous support, non-degeneracy	LP-based	$\mathcal{O}(\log T \log \log T)$	Delay
(Bray, 2022)	Bounded, continuous support, non-degeneracy	LP-based	$\mathcal{O}(\log T)$	Delay
(Chen et al., 2024)	Bounded, finite support, non-degeneracy	LP-based	$\mathcal{O}(1)$	Delay
(Li et al., 2024)	Bounded, finite support, non-degeneracy	LP-based	$\mathcal{O}(1)$	Delay
(Xu et al., 2024)	Bounded, continuous support, non-degeneracy	LP-based	$\mathcal{O}(\log(T/f))$	Delay
<b>This paper</b>	<b>Bounded, continuous support, non-degeneracy</b>	<b>LP-based &amp; First-order</b>	<b><math>\mathcal{O}(\log(T/f) + \sqrt{f})</math></b>	<b>No Delay</b>
(Li et al., 2020)	Bounded	First-order	$\mathcal{O}(\sqrt{T})$	No Delay
(Balseiro et al., 2022a)	Bounded	First-order	$\mathcal{O}(\sqrt{T})$	No Delay
(Gao et al., 2023)	Bounded	First-order	$\mathcal{O}(\sqrt{T})$	No Delay
(Sun et al., 2020)	Bounded, finite support, non-degeneracy	First-order	$\mathcal{O}(T^{3/8})$	No Delay
(Gao et al., 2024)	Bounded, continuous support, non-degeneracy	First-order	$\mathcal{O}(T^{1/3})$	No Delay

**Paper organization.** The rest of the paper is organized as follows. Section 2 introduces the problem formulation and assumptions. Section 3 proposes our algorithms and main theoretical guarantee: a  $\mathcal{O}(\log(T/f) + \sqrt{f})$  regret bound. Section 4 presents experiments to validate our theory.

## 2. Problem Setup

We use  $\|\cdot\|$  to denote the Euclidean norm and  $\langle \cdot \rangle$  to denote Euclidean inner product. Bold letters  $\mathbf{A}$  and  $\mathbf{a}$  denote matrices and vectors, respectively. Subscript  $(\cdot)_{it}$  denotes the index for resource type  $i$  at time  $t$ . The notation  $(\cdot)^+ = \max\{\cdot, 0\}$  denotes the element-wise positive part function, and  $\mathbb{I}(\cdot)$  denotes the 0-1 indicator function.

### 2.1. OLP Formulation

We study an online resource allocation problem over the time horizon  $T$  under a *stochastic input model*. Initially, we have an inventory vector  $\mathbf{b} \in \mathbb{R}^m$ , representing  $m$  resource types. At each time step  $t$ , a customer arrives with a request sampled i.i.d. as  $(r_t, \mathbf{a}_t) \sim \mathcal{P}$ , where  $\mathbf{r} = (r_1, \dots, r_T)^\top \in \mathbb{R}^T$  is the offered payment (bid),  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_T) \in \mathbb{R}^{m \times T}$  is the matrix of customers' resource demands, and  $\mathcal{P}$  is a fixed, unknown distribution. We must decide whether to accept or reject each request, represented by the decision variables  $\mathbf{x} = (x_1, \dots, x_T) \in \{0, 1\}^T$ . The goal is to maximize the cumulative reward. This problem can be formulated as the following OLP, referred to as the *primal linear program (PLP)*:

$$\max_{0 \leq \mathbf{x} \leq \mathbf{1}} \langle \mathbf{r}, \mathbf{x} \rangle \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}. \quad (\text{PLP})$$

The dual problem of (PLP) is given by

$$\min_{(\mathbf{p}, \mathbf{y}) \geq \mathbf{0}} \langle \mathbf{b}, \mathbf{p} \rangle + \langle \mathbf{1}, \mathbf{y} \rangle \quad \text{s.t.} \quad \mathbf{A}^\top \mathbf{p} + \mathbf{y} \geq \mathbf{r} \quad (\text{DLP})$$

where  $\mathbf{p}$  is the vector of *dual prices*. Let  $\mathbf{d} = \mathbf{b}/T \in \mathbb{R}^m$  denote the initial average resource. As demonstrated by (Li

et al., 2020), (DLP) can be written as

$$\min_{\mathbf{p} \geq \mathbf{0}} f_T(\mathbf{p}) := \frac{1}{T} \sum_{t=1}^T [\langle \mathbf{d}, \mathbf{p} \rangle + (r_t - \langle \mathbf{a}_t, \mathbf{p} \rangle)^+] \quad (1)$$

This formulation can be written as a  $T$ -sample approximation of the following stochastic program:

$$\min_{\mathbf{p} \geq \mathbf{0}} f(\mathbf{p}) := \mathbb{E}[f_T(\mathbf{p})] = \mathbf{d}^\top \mathbf{p} + \mathbb{E}[(r - \mathbf{a}^\top \mathbf{p})^+] \quad (2)$$

where the expectation is taken with respect to  $(r, \mathbf{a}) \sim \mathcal{P}$ .

We define optimal solutions to the  $T$ -sample approximation problem (1) and stochastic program (2) respectively as

$$\mathbf{p}_T^* = \arg \min_{\mathbf{p} \geq \mathbf{0}} f_T(\mathbf{p}) \quad \text{and} \quad \mathbf{p}^* = \arg \min_{\mathbf{p} \geq \mathbf{0}} f(\mathbf{p}).$$

The decision variable  $x_t^*$  can be established from the complementary slackness condition as

$$x_t^* = \begin{cases} 0, & r_t < \mathbf{a}_t^\top \mathbf{p}_T^*, \\ 1, & r_t > \mathbf{a}_t^\top \mathbf{p}_T^* \end{cases} \quad (3)$$

and  $x_t^* \in [0, 1]$  if  $r_t = \mathbf{a}_t^\top \mathbf{p}_T^*$ .

This connection between primal and dual solutions inspires OLP algorithms that make decisions based on dual prices:

$$x_t = \mathbb{I}(r_t \geq \mathbf{a}_t^\top \mathbf{p}_t). \quad (4)$$

### 2.2. Algorithms for OLP

We summarize two main dual-based OLP algorithms.

**LP-based Method.** The LP-based method (Li & Ye, 2022) calculates dual prices by re-solving the online LP problem at each time step. Specifically, define  $d_{it} = b_{it}/(T - t)$  as the average remaining resource for type  $i$  at time  $t$ . The resulting optimization problem can be viewed as a  $t$ -sample approximation to the stochastic program specified by  $\mathbf{d}_t = (d_{1t}, \dots, d_{mt})^\top$  as

$$\min_{\mathbf{p} \geq \mathbf{0}} f_{\mathbf{d}_t}(\mathbf{p}) := \mathbf{d}_t^\top \mathbf{p} + \mathbb{E}[(r - \mathbf{a}^\top \mathbf{p})^+] \quad (5)$$

By updating  $d_t$ , this method incorporates past decisions when computing  $p_t$ . If resources are over-utilized in earlier periods, the supply decreases, prompting us to raise the dual price and be more selective with future orders. We outline this method in Algorithm 3 of Appendix A.2.

**First-order Method.** The first-order method (Li et al., 2020) calculates dual prices via online subgradient updates. We maintain a static average resource  $d = b/T$ , compute  $x_t$  as per (4), and update the dual price as

$$p_{t+1} = (p_t - \alpha_t(d - a_t x_t))^+$$

where the subgradient term evaluated at  $p_t$  is

$$d - a_t x_t \in \partial_{p=p_t}(d^\top p + (r_t - a_t^\top p)^+).$$

This process can be interpreted as a projected stochastic subgradient method for solving (1). It reduces computational cost by requiring only a single pass through the data and eliminates solving LPs explicitly. A restart strategy improves the first-order method to have  $\mathcal{O}(T^{1/3})$  regret (Gao et al., 2024). Algorithm 4 and Algorithm 5 in Appendix A.3 summarize this method.

### 2.3. Performance Metrics

We aim to design algorithms that optimize a bi-objective performance measure involving regret and constraint violation. The regret measures the difference between the objective value of the algorithm's output and that of the true optimal solution, while the constraint violation measures the degree to which the algorithm's output fails to meet the given constraints. We denote the offline optimal solution to (PLP) by  $x^* = (x_1^*, \dots, x_T^*)$ , and the online algorithm output by  $x = (x_1, \dots, x_T)$ . Then, we define the regret  $r(x)$  and resource violation  $v(x)$  as

$$r(x) := \langle r, x^* \rangle - \langle r, x \rangle, \quad (6)$$

$$v(x) := \|(Ax - b)^+\|. \quad (7)$$

Therefore, we define the following bi-objectives for evaluating an algorithm's worst-case performance:

$$\Delta_T = \sup_{p \in \Xi} \mathbb{E}_p[r(x) + v(x)] \quad (8)$$

where  $\Xi$  denotes a family of distributions satisfying regularity assumptions specified later.

This metric is commonly used for first-order OLP algorithms (Gao et al., 2023; Li et al., 2020) and is also aligned with the literature on online convex optimization with constraints (Mahdavi et al., 2012; Yu et al., 2017). Integrating  $r(x)$  and  $v(x)$  into a single performance measure promotes balanced resource consumption over the decision horizon. In Section 3, we derive a decomposition of (8) that unifies the regret analysis for all OLP methods.

### 2.4. Assumptions and Auxiliary Results

We adopt the following assumptions regarding the stochastic inputs. These assumptions are standard in the online learning literature (Li & Ye, 2022; Jiang et al., 2022; Xu et al., 2024). In particular, we require the input data to be bounded, follow a linear growth, and be non-degenerate.

**Assumption 2.1** (Boundedness). We assume

- (a) The order inputs  $\{(r_t, a_t)\}_{t=1}^T$  are generated i.i.d from an unknown distribution  $\mathcal{P}$ .
- (b) There exist constants  $\bar{r}, \bar{a} > 0$  such that  $|r_t| \leq \bar{r}$  and  $\|a_t\|_\infty \leq \bar{a}$  almost surely for  $t = 1, \dots, T$ .
- (c) The average resource capacity  $d = b/T$  satisfies  $d_i \in [\underline{d}, \bar{d}]$  for some  $\underline{d} > 0$  for any  $i = 1, \dots, m$ .

In this assumption, (a) states that  $\{(r_t, a_t)\}_{t=1}^T$  are independent of each other, but we allow dependencies between their components. Part (b) introduces the bounds  $\bar{r}, \bar{a}$  solely for analytical purposes. This is a minimal requirement on  $(r_t, a_t)$  compared to previous work (Agrawal et al., 2014; Li & Ye, 2022; Xu et al., 2024). Part (c) requires the average resource to grow linearly with  $T$ , ensuring that a constant fraction of  $x_t$  values can be set to 1. Consequently, the number of fulfillable orders is proportional to  $T$ , facilitating a stable service level over time.

**Assumption 2.2** (Uniform Non-degeneracy). We assume

- (a) The second-moment matrix  $\mathbb{E}[aa^\top]$  is positive definite with minimum eigenvalue  $\lambda$ .
- (b) There are constants  $\mu, \nu$  such that for any  $(r, a) \sim \mathcal{P}$ ,

$$\begin{aligned} & \nu |a^\top(p - p^*)| \\ & \leq |\mathbb{P}(r \geq a^\top p \mid a) - \mathbb{P}(r \geq a^\top p^* \mid a)| \\ & \leq \mu |a^\top(p - p^*)| \end{aligned}$$

holds for all

$$p \in \mathcal{V}_p := \left\{ p \in \mathbb{R}^m : p \geq 0, \|p\| \leq \frac{\bar{r}}{\underline{d}} \right\}$$

and  $d \in \mathcal{V}_d := [\underline{d}, \bar{d}]^m$ .

- (c) The optimal  $p^*$  satisfies  $p_i^* = 0$  if and only if  $d_i - \mathbb{E}[a_i \mathbb{I}(r > a^\top p^*)] > 0$  for all  $d \in \mathcal{V}_d$  and  $i \in [m]$ .

In this assumption, part (b) ensures that the cumulative distribution of the reward given the resource consumption request  $r|a$  is continuous and exhibits a stable growth rate. Part (c) requires strict complementarity for the optimal solutions of the stochastic program in (2), which is a non-degeneracy condition for both the primal and dual LPs.



According to stochastic program (2), we define the binding and non-binding index sets as:

$$\begin{aligned} I_B &= \{i : d_i - \mathbb{E}[a_i \mathbb{I}(r > \mathbf{a}^\top \mathbf{p}^*)] = 0\}, \\ I_N &= \{i : d_i - \mathbb{E}[a_i \mathbb{I}(r > \mathbf{a}^\top \mathbf{p}^*)] > 0\}. \end{aligned} \quad (9)$$

By Assumption 2.2(c), these sets are complements, as  $I_B \cap I_N = \emptyset$  and  $I_B \cup I_N = \{1, \dots, m\}$ .

### 3. Parallel Multi-Phase OLP Algorithm

We now present our algorithm for parallel multi-phase online learning and decision-making. Our approach is motivated by the challenges of existing methods. Specifically, the LP-based method has a strong worst-case regret bound of  $\mathcal{O}(\log T)$ , but its high computational costs lead to decision-making delays. Meanwhile, the first-order method updates decisions efficiently but suffers a high regret of  $\mathcal{O}(\sqrt{T})$ . By combining the strengths of these two methods, our new framework balances decision-making quality and computational efficiency.

#### 3.1. Algorithm Design

We establish our framework with the following design:

1. We maintain two parallel paths of online learning and online decision-making. The online learning results are periodically sent to the decision-making path as a re-start point to guide subsequent updates.
2. In the online learning path, we employ a streamlined LP-based method, which only re-solves the updated OLP problem according to a predefined frequency, reducing the computational overhead.
3. In the online decision-making path, we apply the first-order method during the initial and final batches. The learning rate is optimally tuned for these two intervals.

Figure 1 illustrates our framework, which generates two parallel sequences of dual prices:  $\{\mathbf{p}_t^D\}_{t=1}^T$  from the first-order method and  $\{\mathbf{p}_t^L\}_{t=1}^T$  from the LP-based method. The algorithm proceeds in batches of length  $f$ . In the first batch, it uses the first-order method to iteratively update the dual prices, thereby guiding decision-making. At the end of this batch, it applies the LP-based method to obtain a refined dual price and passes it to the decision-making path. The LP re-solving occurs only once per batch, making the approach computationally efficient. Formally, we re-solve the OLP at every time  $t$  satisfying  $t \leq kf$  and  $t \bmod f = 0$ , where  $k = \lfloor T/f \rfloor$  is the number of batches.

During intermediate batches, decisions are made based on the most recent dual prices computed from the LP at the end

#### Algorithm 1 Parallel Multi-Phase OLP Algorithm

**Input:** total resource  $\mathbf{b}$ , time horizon  $T$ , average resource  $\mathbf{d} = \mathbf{b}/T$ , initial dual price  $\mathbf{p}_1 = \mathbf{0}$ , re-solving frequency  $f$ , and number of batches  $k = \lfloor T/f \rfloor$

**for**  $t = 1$  to  $T$  **do**

Observe  $(r_t, \mathbf{a}_t)$  and make decision  $x_t$  as rule (4)

Update constraint for  $i = 1, \dots, m$ :

remaining resource  $b_{it} = b_{i,t-1} - a_{it}x_t$

average remaining resource  $d_{it} = \frac{b_{it}}{T-t}$

**if**  $t \leq f$  or  $t \geq kf$  **then**

Update learning rate  $\alpha_t = \begin{cases} \mathcal{O}(1/f^{1/2}) & t \leq f \\ \mathcal{O}(1/f^{2/3}) & t \geq kf \end{cases}$

Compute subgradient and update dual price  $\mathbf{p}_{t+1}$ :

$\mathbf{p}_{t+1} = (\mathbf{p}_t - \alpha_t(\mathbf{d} - \mathbf{a}_t x_t))^+$

**end**

**if**  $t \bmod f = 0$  and constraints are not violated **then**

Solve OLP and update dual price  $\mathbf{p}_{t+1}$ :

$\mathbf{p}_{t+1} = \underset{\mathbf{p} \geq \mathbf{0}}{\operatorname{argmin}} \mathbf{d}_t^\top \mathbf{p} + \frac{1}{t} \sum_{j=1}^t (r_j - \mathbf{a}_j^\top \mathbf{p})^+$

**end**

**else**

Update dual price  $\mathbf{p}_{t+1}$  to be the most recent solution  $\mathbf{p}_{t+1} = \mathbf{p}_t$ .

**end**

**end**

of the previous batch. The algorithm restarts the subgradient updates only in the final batch to guide the remaining decisions. Algorithm 1 summarizes this approach.

Algorithm 1 balances online learning and efficient decision-making. It is responsive to dynamic environments, adapts to the latest information, and reduces computational cost by periodic re-solving. We thus establish a “wait-less” online decision-making framework where each customer’s order is processed immediately without delays from earlier requests or large-scale LP computations.

#### 3.2. Algorithm Analysis

We decompose the performance metric of Algorithm 1 into three key components. All proofs (including essential properties of the dual price  $\mathbf{p}_t$ ) are in Appendix B and C.

**Theorem 3.1** (Performance Metric). *Under Assumptions 2.1 and 2.2, the performance  $\Delta_T$  of Algorithm 1 satisfies*

$$\begin{aligned} \Delta_T &\leq \mu \bar{a}^2 \sum_{t=1}^T \mathbb{E} [\|\mathbf{p}_t - \mathbf{p}^*\|^2] + \mathbb{E} [\|(\mathbf{A}\mathbf{x} - \mathbf{b})^+\|] \\ &\quad + \|\mathbf{p}^*\| \cdot \mathbb{E} [\|(\mathbf{b} - \mathbf{A}\mathbf{x})^{B+}\|]. \end{aligned} \quad (10)$$

where  $(\cdot)^{B+}$  indicates the projection of binding terms onto the positive orthant.

Based on the definition of  $\Delta_T$  (8), we derive the perfor-

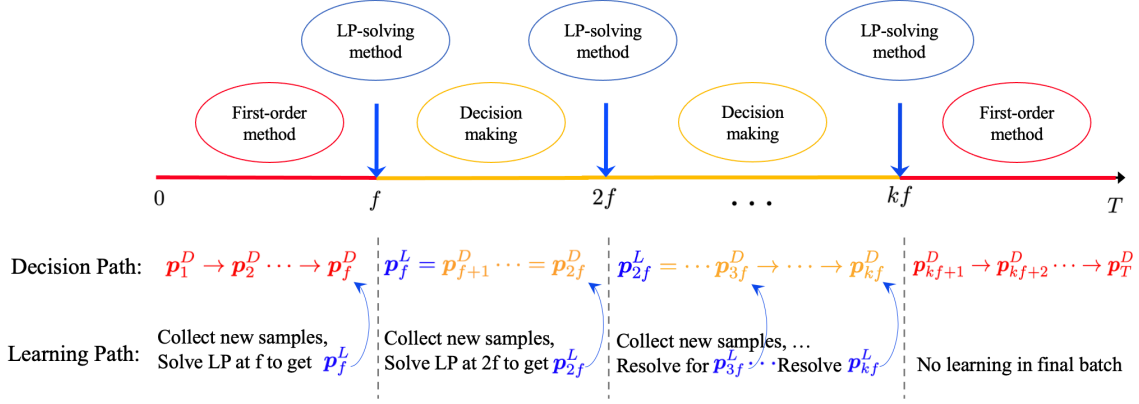


Figure 1: Algorithm 1 illustration of parallel paths and the interactions between online learning and decision-making. Decisions are generated based on 1) the LP-based method (blue) with frequency  $f$ , 2) the first-order method (red) for the initial and final phases (with a warm start), and 3) employing the latest dual price (yellow) during intermediate phases.

mance metric as the sum of three key components: dual convergence, remaining resources, and binding constraint violation. This bound demonstrates that our algorithm encourages 1) smooth and balanced resource utilization, and 2) full resource consumption by the end of the time horizon. Building on Theorem 3.1, we provide a “spectrum theorem” for the algorithm’s performance.

**Theorem 3.2.** *Under Assumptions 2.1 and 2.2, the worst-case performance  $\Delta_T$  of Algorithm 1 is bounded by*

$$\Delta_T \in \mathcal{O}\left(\log\left(\frac{T}{f}\right) + \sqrt{f}\right) \quad (11)$$

where  $T/f$  represents the number of re-solving batches and  $f$  is the length of each batch.

**Remark 3.3** (Spectrum Theorem). We elucidate the trade-offs in total regret induced by using LP-based methods with  $\mathcal{O}(\log T)$  regret and first-order methods with  $\mathcal{O}(\sqrt{T})$  regret. Our algorithm introduces a re-solving frequency  $f \in [1, T]$  and achieves  $\mathcal{O}(\log(T/f) + \sqrt{f})$  regret, which recovers previous results with extreme cases of  $f = 1$  and  $f = T$ . Specifically, the first-order method gives us the regret of  $\sqrt{f}$  for the first batch and  $f^{1/3}$  for the last batch, and the LP-based method contributes the regret of  $\log(T/f)$ .

**Remark 3.4** (Warm Start). We can obtain a tighter regret bound of  $\mathcal{O}(\log(T/f) + f^{1/3})$  given a warm start of the initial dual price satisfying  $\|\mathbf{p}_0 - \mathbf{p}^*\| \leq f^{-1/3}$  or if we use the LP-based method at each time step for the first batch.

**Remark 3.5** (Learning Rate Selection). We select the best learning rate to minimize the regret upper bound consisting of (6) and (7). For the first batch, the regret grows linearly with  $\alpha_t$  while the constraint violation is inversely proportional to  $\alpha_t$ . To achieve the tightest bound, we balance this trade-off by selecting the learning rate that minimizes the overall expression, which yields the optimal choice  $\alpha_t = 1/\sqrt{f}$  as derived in Theorem B.7. A similar anal-

ysis for the final batch leads to the choice  $\alpha_t = 1/f^{2/3}$  as shown in Theorem B.9.

**Technical Intuitions.** Theorem 3.1 decomposes the performance metrics to make the regret analysis compatible between the LP-based and first-order methods.

- **LP-based Method:** As Theorem 3.1 suggests, achieving small regret requires  $\mathbf{A}\mathbf{x} - \mathbf{b}$  to be close to zero. To enforce this, we impose a stronger condition to track the average remaining resource—if  $d_t$  exceeds the allowed limit, we manually set the dual prices to zero to accept all subsequent orders. This approach enables us to eliminate the “stopping time” argument and make the analysis compatible with our new framework.
- **First-order Method:** The regret improvement comes from the final batch since it starts with the LP-derived learning result  $\mathbf{p}_{kf}$ , which lies within a  $\mathcal{O}(1/\sqrt{kf})$ -sized neighborhood around  $\mathbf{p}^*$ . Using the subgradient method, we bound the three components in Theorem 3.1 in terms of  $\|\mathbf{p}_t - \mathbf{p}^*\|$  and express them in terms of  $f, k, T$ , and  $\alpha_t$ . Optimizing the learning rate to  $\alpha_t = f^{-2/3}$  yields an improved regret of  $\mathcal{O}(f^{1/3})$  in the final batch.

### 3.3. Algorithm Extension

We propose an enhanced version of Algorithm 1 that employs the first-order method for decision-making between two consecutive LP resolves. Rather than making decisions solely with the most recently solved dual price from the learning path, Algorithm 2 treats it as a new starting point for each re-solving interval and adopts a smaller step size for subgradient updates to avoid deviating too far from the LP-guided solutions. Figure 2 illustrates this framework. Specifically, for each interval  $[jf, (j+1)f]$ , we attain the new start  $\mathbf{p}_{jf}$  from the learning path at time  $jf$  and con-

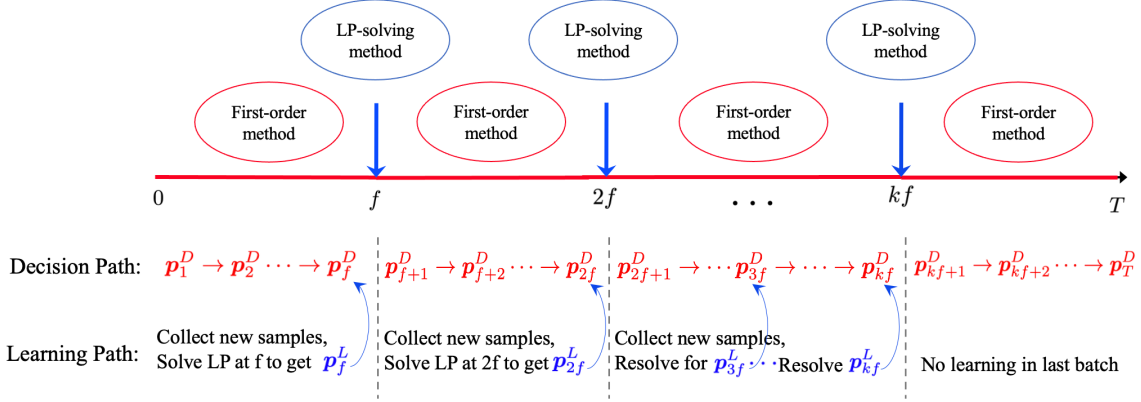


Figure 2: Algorithm 2 illustration of parallel paths with multi-time restart. Decisions are generated based on 1) the LP-based method (blue) with frequency  $f$  and 2) the first-order method (red) during re-solving intervals with a warm start.

#### Algorithm 2 Enhanced Multi-Start OLP Algorithm

**Input:** total resource  $\mathbf{b}$ , horizon  $T$ , average resource  $\mathbf{d} = \mathbf{b}/T$ , initial dual  $\mathbf{p}_1 = \mathbf{0}$ , re-solving frequency  $f$

**for**  $t = 1$  to  $T$  **do**

    Observe  $(\mathbf{r}_t, \mathbf{a}_t)$  and make decision  $\mathbf{x}_t$  as rule (4)

    Update constraint for  $i = 1, \dots, m$ :

        remaining resource constraint  $b_{it} = b_{i,t-1} - a_{it}x_t$

        average resource capacity  $d_{it} = \frac{b_{it}}{T-t}$

**if**  $t \bmod f \neq 0$  **then**

        Update learning rate  $\alpha_t = \mathcal{O}(1/t)$

        Compute subgradient and update dual price  $\mathbf{p}_{t+1}$ :

$\mathbf{p}_{t+1} = (\mathbf{p}_t - \alpha_t(\mathbf{d} - \mathbf{a}_t \mathbf{x}_t))^+$

**end**

**if**  $t \bmod f = 0$  and constraints are not violated **then**

        Solve OLP and update dual price  $\mathbf{p}_{t+1}$ :

$\mathbf{p}_{t+1} = \underset{\mathbf{p} \geq \mathbf{0}}{\operatorname{argmin}} \mathbf{d}_t^\top \mathbf{p} + \frac{1}{t} \sum_{j=1}^t (\mathbf{r}_j - \mathbf{a}_j^\top \mathbf{p})^+$

**end**

**end**

tinue to fine-tune the dual price  $\{\mathbf{p}_t\}_{t=jf}^{(j+1)f}$  using the first-order method with a learning rate  $\alpha_t = \mathcal{O}(1/t)$ . Algorithm 2 summarizes this approach.

As we illustrate in our experiments (Section 4), the incorporation of an intermediate first-order method improves the algorithm's stability and ensures smooth resource consumption. Therefore, the multi-restart mechanism results in better performance during the final batch and improves the algorithm's total performance.

#### 3.4. Algorithm Application

Our motivation in designing Algorithm 1 is to effectively balance computational efficiency and decision optimality.

Building on the Spectrum Theorem 3.2, this section aims to translate our theoretical results into practical applications.

We formulate a new optimization problem to find the optimal re-solving frequency that minimizes the regret in Theorem 3.2, taking into account computational resource capacities. Denote  $c_1(\cdot), c_2(\cdot)$  as computational cost functions for LP-based and first-order methods, and  $R$  as the total computational resource capacity. The optimal value of  $f$  is determined by solving the following optimization problem:

$$\begin{aligned} \min_{f \in \{1, \dots, T\}} \quad & \log\left(\frac{T}{f}\right) + \sqrt{f} \\ \text{s.t.} \quad & c_1(k) + 2c_2(f) \leq R. \end{aligned} \quad (12)$$

Specifically, if we use the interior-point method or the simplex method as the LP solver, the computational cost is  $m^2(m+t)$  for any time  $t$ . The first-order method updates gradients in constant time. The following proposition provides a concrete example in practice.

**Proposition 3.6** (Optimal Re-solving Frequency). *Given a fixed computation resource capacity  $R$ , if we use the interior-point or simplex method as the LP solver in Algorithm 1, we can instantiate Constraint (12) as*

$$\sum_{j=1}^k (m^2(m+jf)) + 2mf \leq R.$$

This proposition enables users to determine the optimal re-solving frequency that balances regret and computational cost based on available computational resources.

## 4. Numerical Experiments

We conduct extensive experiments to evaluate our algorithm's performance and validate our theoretical results. This section is divided into two parts. In the first part (Section 4.1), we evaluate Algorithms 1 and 2 across different choices of re-solving frequency. In the second part (Section 4.2), we compare our algorithm with LP-based and

Table 3: Algorithms comparison.

$T$	Regret	Algorithm	Compute Time (s)
$10^4$	115.32	$\mathcal{O}(T^{1/2})$ First-Order	0.008
	60.39	$\mathcal{O}(T^{1/3})$ First-Order	0.013
	3.50	LP-based	123.497
	<b>9.12</b>	Algorithm 1	<b>1.8</b>
	<b>5.09</b>	Algorithm 2	<b>1.8</b>
$10^5$	203.20	$\mathcal{O}(T^{1/2})$ First-Order	0.118
	73.07	$\mathcal{O}(T^{1/3})$ First-Order	0.108
	3.88	LP-based	> 3600
	<b>9.41</b>	Algorithm 1	<b>56.9</b>
	<b>6.36</b>	Algorithm 2	<b>56.8</b>
$10^6$	351.91	$\mathcal{O}(T^{1/2})$ First-Order	1.211
	115.39	$\mathcal{O}(T^{1/3})$ First-Order	1.577
	5.50	LP-based	> 100000
	<b>11.65</b>	Algorithm 1	<b>2155.9</b>
	<b>7.09</b>	Algorithm 2	<b>2242.1</b>

first-order methods in terms of regret and running time. All implementations can be found at [GitHub Link](#).

We consider the following distributions:

Input I:  $a_{it} \sim \text{Unif}[0, 2], r_t \sim \text{Unif}[0, 10]$

Input II:  $a_{it} \sim \mathcal{N}(0.5, 1), r_t \sim \mathcal{N}(0.5m, m)$

Learning rates are selected as specified in Section 3:

$$\text{Algorithm 1: } \alpha_t = \begin{cases} \mathcal{O}(1/f^{1/2}) & t \leq f \\ \mathcal{O}(1/f^{2/3}) & t \geq kf \end{cases}$$

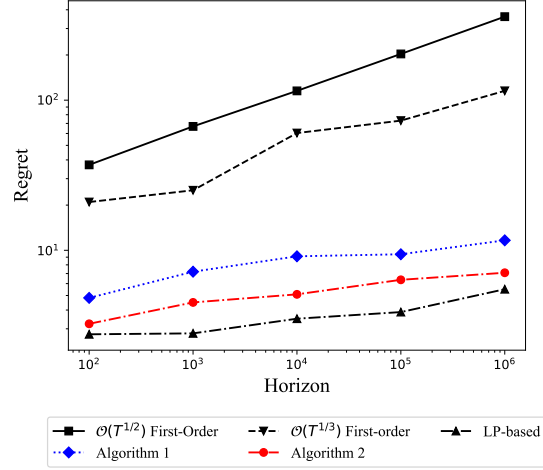
$$\text{Algorithm 2: } \alpha_t = \mathcal{O}(1/t)$$

#### 4.1. Regret under Varying Re-solving Frequencies

We choose  $m = 1$  and generate the sequence  $\{(r_t, \mathbf{a}_t)\}_{t=1}^T$  randomly from the uniform (Input I) and normal (Input II) distributions, and include more complex distributions in Appendix E.1. The time horizon  $T$  spans evenly over  $[10^2, 10^6]$ , the initial average resource is sampled as  $d_i \sim \text{Uniform}[1/3, 2/3]$ , and  $f \in \{T^{1/3}, T^{1/2}, T^{2/3}\}$  representing high, medium, and low re-solving frequencies. We report the average result over 100 trials for each experiment. We use the classic first-order method with  $\mathcal{O}(T^{1/2})$  regret (Li et al., 2020) as a baseline. Results are summarized in Table 2 and Figure 3 on a logarithmic scale.

We analyze the absolute value of regret and its growth rate over time. We observe that regret decreases as the re-solving frequency increases. This trend holds consistently across both algorithms and input types. In addition, while regret accumulates over longer time horizons, the increasing rate remains stable for algorithms employing higher re-solving frequencies. These findings are consistent with the guarantees of Theorem 3.2, as more frequent updates enable better adaptation to dynamic environments.

Figure 4: Regret for various algorithms.



#### 4.2. Comparative Analysis with Baseline Methods

We next compare our algorithm’s regret and computation time with a classic LP-based (Li & Ye, 2022) and two first-order methods with  $\mathcal{O}(T^{1/2})$  (Li et al., 2020) and  $\mathcal{O}(T^{1/3})$  regrets (Gao et al., 2024). We generate  $\{(r_t, \mathbf{a}_t)\}_{t=1}^T$  from a uniform distribution (Input I). We set the resource types to  $m = 5$ , time  $T$  to range evenly over  $[10^2, 10^6]$ , average resource  $d_i \sim \text{Uniform}[1/3, 2/3]$ , and re-solving frequency to  $f = T^{1/3}$  from Section 4.1. Each result is averaged over 100 trial runs. We summarize the findings in Table 3 and Figure 4, and we include more comparisons with an infrequent re-solving method (Li et al., 2024) in Appendix E.2. We observe the effectiveness of Algorithm 1 and 2:

1. Our algorithms exhibit strong performance in terms of decision optimality. They achieve over a 20-fold and 10-fold improvement in regret compared to the  $\mathcal{O}(T^{1/2})$  first-order method and the  $\mathcal{O}(T^{1/3})$  first-order method respectively. These numerical results also corroborate the theoretical bounds in Theorem 3.2.
2. Our algorithms are computationally efficient. Their run-times exhibit over 100-fold improvement compared to the LP-based method with only a minimal increase (less than 2-fold) in regret.

Therefore, we achieve a balance between effective decision-making and efficient computation. Our algorithms demonstrate better regret than the first-order method and obtain substantial computational speed-ups compared to the LP-based method. They exhibit strong scalability and adaptability across various re-solving frequencies and stochastic input models, consistently delivering superior performance as the problem size grows.



Table 2: Regret of algorithms with various re-solving frequencies.

	$T$	First-Order	Low freq	Mid freq	High freq
Input I	$10^3$	12.13	7.76	5.77	<b>4.86</b>
	$10^4$	38.50	10.96	7.55	<b>5.67</b>
	$10^5$	122.44	23.90	9.92	<b>8.36</b>
	$10^6$	404.59	56.70	21.90	<b>8.99</b>
Input II	$10^3$	11.44	6.28	4.86	<b>3.95</b>
	$10^4$	36.50	10.21	7.34	<b>3.81</b>
	$10^5$	115.57	14.61	11.78	<b>4.66</b>
	$10^6$	365.99	35.20	15.68	<b>6.26</b>

	$T$	First-Order	Low freq	Mid freq	High freq
Input I	$10^3$	12.13	6.78	5.20	<b>4.50</b>
	$10^4$	38.50	10.37	8.03	<b>5.99</b>
	$10^5$	122.44	22.33	11.57	<b>6.36</b>
	$10^6$	404.59	48.21	22.44	<b>7.09</b>
Input II	$10^3$	11.44	3.20	2.56	<b>1.75</b>
	$10^4$	36.50	5.48	4.30	<b>2.52</b>
	$10^5$	115.57	12.35	4.48	<b>3.86</b>
	$10^6$	365.99	30.48	13.20	<b>4.77</b>

(a) Algorithm 1
(b) Algorithm 2

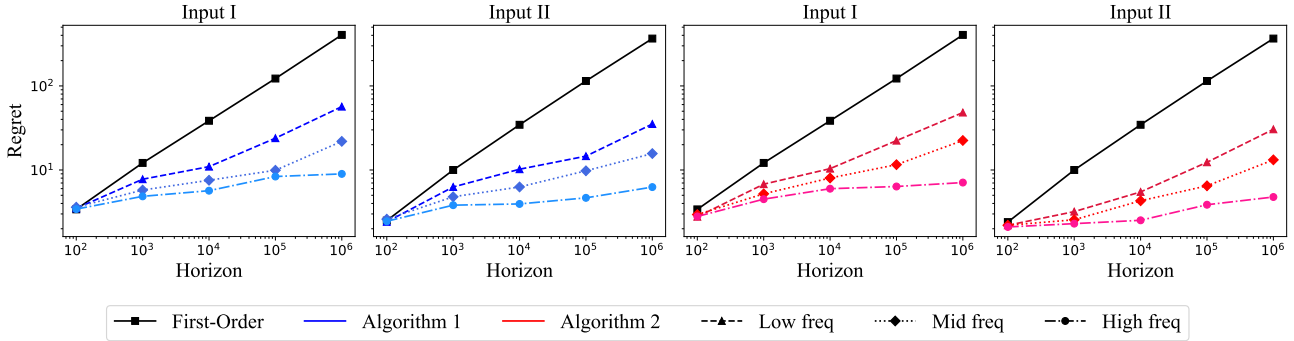


Figure 3: Evaluations of Algorithm 1 and 2 across various horizons, re-solving frequencies, and stochastic inputs, validating the positive relationship between regret and frequency stated in Theorem 3.2.

## 5. Conclusion and Discussion

This paper presents a new approach to online linear programming for dynamic resource allocation, addressing the inherent trade-offs between decision-making optimality and computational efficiency. Recognizing the limitations of existing methods, we propose a parallel framework that decouples online learning and decision-making into independent yet complementary processes. By integrating LP-based and first-order methods, our framework effectively balances total regret and computational cost.

We establish rigorous theoretical guarantees, proving that our algorithm achieves a worst-case regret bound of  $\mathcal{O}(\log(T/f) + \sqrt{f})$  under continuous support. This result highlights our method’s ability to interpolate between LP-based and first-order methods based on computational capability. Furthermore, extensive experiments validate the effectiveness of our approach, demonstrating improvements in both regret minimization and runtime efficiency over competitive baselines.

Beyond these contributions, our work opens avenues for further research in adaptive online decision-making. Future directions include refining the re-solving frequency based on real-time computational constraints and extending our

framework to broader classes of online optimization problems. Overall, our results underscore the potential of hybrid algorithms in high-dimensional, large-scale environments, offering practical insights for applications in operations research, machine learning, and beyond.

### Discussion of Unknown Horizon

In this paper, we consider decision-making under a finite horizon, since the average resource Assumption 2.1(c) widely used in OLP literature relies on a fixed horizon to be well-defined. When the horizon is unknown, this assumption becomes ill-posed, and prior work (Balseiro et al., 2023) shows that it may not be possible to achieve sub-linear regret. Addressing these challenges would likely require first adapting LP-based and first-order OLP methods to uncertain horizons before tackling a hybrid approach.

In practice, it is often possible to make a prior prediction of the horizon using data-driven approaches or based on the resources available. For example, in online advertising, the number of customers can be predicted from historical statistics. Besides, in retail or event-based sales, the selling horizon may be externally decided by upper-level decision-makers. These practical applications suggest a reasonable modeling choice for finite time horizons.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## Acknowledgments

This work was supported in part by NSF grant CCF-2338226.

## References

- Agrawal, S., Wang, Z., and Ye, Y. A dynamic near-optimal algorithm for online linear programming. *Operations Research*, 62(4):876–890, 2014.
- Balseiro, S., Kroer, C., and Kumar, R. Online resource allocation under horizon uncertainty. In *Abstract Proceedings of the 2023 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, pp. 63–64, 2023.
- Balseiro, S. R., Lu, H., and Mirrokni, V. The best of many worlds: Dual mirror descent for online allocation problems. *Operations Research*, 2022a.
- Balseiro, S. R., Lu, H., Mirrokni, V., and Sivan, B. From online optimization to PID controllers: Mirror descent with momentum. *arXiv preprint arXiv:2202.06152*, 2022b.
- Bray, R. L. Logarithmic regret in multisectionary and online linear programming problems with continuous valuations. *arXiv preprint arXiv:1912.08917*, 2022.
- Chen, G., Li, X., and Ye, Y. An improved analysis of lp-based control for revenue management. *Operations Research*, 72(3):1124–1138, 2024.
- Devanur, N. R., Jain, K., Sivan, B., and Wilkens, C. A. Near optimal online algorithms and fast approximation algorithms for resource allocation problems. *Journal of the ACM (JACM)*, 66(1):1–41, 2019.
- Gao, W., Ge, D., Sun, C., and Ye, Y. Solving linear programs with fast online learning algorithms. In *International Conference on Machine Learning*, pp. 10649–10675. PMLR, 2023.
- Gao, W., Sun, C., Xue, C., and Ye, Y. Decoupling learning and decision-making: Breaking the  $\mathcal{O}(\sqrt{T})$  barrier in online resource allocation with first-order methods. In *International Conference on Machine Learning*, pp. 14859–14883. PMLR, 2024.
- Goel, G. and Mehta, A. Online budgeted matching in random input models with applications to adwords. In *SODA*, volume 8, pp. 982–991, 2008.
- Golrezaei, N. and Yao, E. Upfront commitment in online resource allocation with patient customers. *arXiv preprint arXiv:2108.03517*, 2021.
- Gupta, A. and Molinaro, M. How the experts algorithm can help solve lps online. *Mathematics of Operations Research*, 41(4):1404–1431, 2016.
- Jiang, J., Ma, W., and Zhang, J. Degeneracy is ok: Logarithmic regret for network revenue management with indiscrete distributions. *arXiv preprint arXiv:2210.07996*, 2022.
- Li, G., Wang, Z., and Zhang, J. Infrequent resolving algorithm for online linear programming. *arXiv preprint arXiv:2408.00465*, 2024.
- Li, X. and Ye, Y. Online linear programming: Dual convergence, new algorithms, and regret bounds. *Operations Research*, 70(5):2948–2966, 2022.
- Li, X., Sun, C., and Ye, Y. Simple and fast algorithm for binary integer and online linear programming. *Advances in Neural Information Processing Systems*, 33:9412–9421, 2020.
- Ma, W., Cao, Y., Tsang, D. H., and Xia, D. Optimal regularized online allocation by adaptive re-solving. *Operations Research*, 2024.
- Mahdavi, M., Jin, R., and Yang, T. Trading regret for efficiency: Online convex optimization with long term constraints. *Journal of Machine Learning Research*, 13(81):2503–2528, 2012.
- Mehta, A., Saberi, A., Vazirani, U., and Vazirani, V. Adwords and generalized online matching. *Journal of the ACM (JACM)*, 54(5):22–es, 2007.
- Sun, R., Wang, X., and Zhou, Z. Near-optimal primal-dual algorithms for quantity-based network revenue management. *arXiv preprint arXiv:2011.06327*, 2020.
- Talluri, K. T., Van Ryzin, G., and Van Ryzin, G. *The theory and practice of revenue management*, volume 1. Springer, 2004.
- Xie, Y., Ma, W., and Xin, L. The benefits of delay to online decision-making. *Available at SSRN 4248326*, 2023.
- Xu, H., Glynn, P. W., and Ye, Y. Online linear programming with batching. *arXiv preprint arXiv:2408.00310*, 2024.

Yu, H., Neely, M., and Wei, X. Online convex optimization with stochastic constraints. *Advances in Neural Information Processing Systems*, 30, 2017.

# Appendix

## Table of Contents

---

<b>A Primary Results and Properties</b>	<b>13</b>
A.1 Preliminary results . . . . .	13
A.2 LP-based analysis . . . . .	13
A.3 First-order Analysis . . . . .	14
<b>B Algorithm Design and Analysis</b>	<b>17</b>
B.1 Regret Analysis . . . . .	17
B.2 Total Performance Analysis . . . . .	18
B.3 Regret for LP-based Method . . . . .	19
B.4 Regret for First-Order Method . . . . .	21
<b>C Main Results</b>	<b>25</b>
C.1 Proof of Theorem 3.1. . . . .	25
C.2 Proof of Theorem 3.2. . . . .	25
<b>D Auxiliary Results</b>	<b>26</b>
D.1 Technical Support for LP-based Analysis . . . . .	26
D.2 Technical Support for First-order Analysis . . . . .	29
<b>E Supplementary Experiments</b>	<b>30</b>
E.1 New Distribution . . . . .	30
E.2 More Comparison . . . . .	30

---

**Structure of the Appendix** We organize the appendix as follows. In **Section A**, we present some primary results for the OLP problem and the foundation analysis for LP-based and first-order methods; **Section B** characterizes the unique properties from our parallel multi-phase structure, and presents the cohesive mathematical framework which decomposes the performance metrics to unify all the OLP methods; **Section C** demonstrates the outline of Algorithm 2, combines all the previous results and proves the main results in our paper; **Section D** provides auxiliary results and technical support for the previous three sections; **Section E** includes supplementary experiments of more general input distributions and comparisons with recent algorithms.

## A. Primary Results and Properties

In this section, we present some primary results for the problem and OLP dual algorithms which will help present our main results. We also include basic properties and convergence of the LP-based and first-order methods. All the analyses are stated under Assumption 2.1 and 2.2.

### A.1. Preliminary results

We propose the following lemma as directed results from Assumptions. Let  $\Xi$  denote the family of distributions satisfying Assumption 2.1 and 2.2. We propose the following lemma.

**Lemma A.1** (Dimension Stability). *Under Assumptions 2.1 and 2.2, there exists a constant  $\delta > 0$  such that*

$$\forall \mathbf{d}_t \in \mathcal{D} := [d_i - \delta, d_i + \delta]^m, \text{ stochastic programs (5) specified by } \mathbf{d}_t \text{ share the same } I_B \text{ and } I_N \text{ sets.}$$

This lemma is a consequence of Lemmas 12 and 13 in (Li & Ye, 2022). The existence of  $\delta$  comes from the continuity of  $f_{\mathbf{d}_t}(\mathbf{p})$ . Thus,  $\delta$  is only associated with stochastic program (5), and it is independent of  $T$ . Note that  $\mathcal{D} \subset \mathcal{V}_{\mathbf{d}}$ . We will analyze our feasible dual solutions derived from the online algorithm with  $\mathbf{d}_t \in \mathcal{D}$ .

**Lemma A.2** (Quadratic Regularity, Proposition 2 in (Li & Ye, 2022)). *Under Assumptions 2.1 and 2.2, for any  $\mathbf{p} \in \mathcal{V}_{\mathbf{p}}$ ,*

$$\begin{aligned} f(\mathbf{p}) &\leq f(\mathbf{p}^*) + \nabla f(\mathbf{p}^*)^\top (\mathbf{p} - \mathbf{p}^*) + \frac{\mu \bar{a}^2}{2} \|\mathbf{p} - \mathbf{p}^*\|^2, \\ f(\mathbf{p}) &\geq f(\mathbf{p}^*) + \nabla f(\mathbf{p}^*)^\top (\mathbf{p} - \mathbf{p}^*) + \frac{\nu \lambda}{2} \|\mathbf{p} - \mathbf{p}^*\|^2. \end{aligned}$$

Moreover,  $\mathbf{p}^*$  is the unique optimal solution to (2).

This lemma establishes a local form of semi-strong convexity and smoothness at  $\mathbf{p}^*$ , which is guaranteed by our assumptions on the distribution  $\mathcal{P} \in \Xi$ . By focusing on these local properties rather than insisting on strong convexity and global smoothness, we relax the classical requirements typically imposed in such settings. In later sections, we will leverage this result to derive our regret bound.

### A.2. LP-based analysis

We state the LP-based method in Algorithm 3 and its properties of boundedness and dual convergence results.

---

#### Algorithm 3 LP-based method

---

**Input:** total resource  $\mathbf{b}$ , time horizon  $T$ , average resource  $\mathbf{d} = \mathbf{b}/T$ , and initial dual price  $\mathbf{p}_1 = \mathbf{0}$ .

**for**  $t = 1$  to  $T$  **do**

    Observe  $(r_t, \mathbf{a}_t)$  and make decision  $x_t$  based on rule (4) if constraints are not violated.

    Update constraint for  $i = 1, \dots, m$ :

        remaining resource constraint  $b_{it} = b_{i,t-1} - a_{it}x_t$

        average resource capacity  $d_{it} = \frac{b_{it}}{T-t}$

    Solve the updated dual problem and obtain dual price  $\mathbf{p}_{t+1}$ :

$$\mathbf{p}_{t+1} = \underset{\mathbf{p} \geq \mathbf{0}}{\operatorname{argmin}} \mathbf{d}_t^\top \mathbf{p} + \frac{1}{t} \sum_{j=1}^t (r_j - \mathbf{a}_j^\top \mathbf{p})^+$$

**end**

---

LP-based method incorporates the past decisions into the optimization of  $\mathbf{p}_t$ . If resources were over-utilized in earlier periods, the remaining supply decreases, prompting the algorithm to raise the dual price and become more selective with future orders. Conversely, if ample resources remain, the future dual price will be lowered, allowing more consumer requests to be accepted. This adaptive mechanism accounts for past actions by adjusting the available resource capacity.

**Lemma A.3** (Boundedness of LP result). *The online dual price  $\mathbf{p}_t$  from Algorithm 3 and the optimal dual price  $\mathbf{p}^*$  of stochastic program (2) are bounded as*

$$\|\mathbf{p}^*\| \leq \frac{\bar{r}}{\underline{d}}, \quad \|\mathbf{p}_t\| \leq \frac{\bar{r}}{\underline{d} - \delta}.$$



**Lemma A.4** (Dual Convergence of LP-based algorithm). *Under Assumptions 2.1 and 2.2,  $\mathbf{p}_t$  represents the online solution from LP-based method, there exists a constant  $C_{lp} > 0$  depending on  $\bar{r}, \bar{a}, \underline{d}, m, \nu$ , and  $\lambda$  such that*

$$\mathbb{E} [\|\mathbf{p}_t - \mathbf{p}^*\|^2] \leq \frac{C_{lp}}{t}. \quad (13)$$

*In addition, the difference between  $\mathbf{p}_t^*$  and  $\mathbf{p}^*$  satisfies*

$$\|\mathbf{p}_t^* - \mathbf{p}^*\|^2 \leq \frac{1}{\nu^2 \lambda^2} \|\mathbf{d}_t - \mathbf{d}\|^2. \quad (14)$$

This lemma establishes that the LP-based online AHDL algorithm (Jiang et al., 2022) produces dual solutions  $\mathbf{p}_t$  that converge to  $\mathbf{p}^*$ . This convergence highlights the stability of the online dual variables during the intermediate stages of decision-making and ensures a warm start for the first-order method in the last re-solving batch. Moreover, we bound the distance between  $\mathbf{p}_t^*$  and  $\mathbf{p}^*$  by relating their respective average resource capacities,  $\mathbf{d}_t$  and  $\mathbf{d}$ , in the associated stochastic programs (Li & Ye, 2022). These results provide the foundation for analyzing  $\|\mathbf{p}_t - \mathbf{p}^*\|$  in our final result.

**Proof of Lemma A.3.** By the optimality of  $\mathbf{p}^*$  and boundedness of  $\mathbf{d}$ , we have

$$\underline{d} \|\mathbf{p}^*\|_1 \leq \mathbf{d}^\top \mathbf{p}^* \leq \mathbb{E}[r] \leq \bar{r}.$$

This holds because if otherwise,  $\mathbf{p}^*$  can not be the optimal solution due to  $f(\mathbf{p}^*) > f(\mathbf{0})$ . Given the non-negativeness of  $\mathbf{p}^*$  and  $\mathbf{p}^* \in \mathcal{V}_{\mathbf{d}}$ , we obtain  $\|\mathbf{p}^*\| \leq \|\mathbf{p}^*\|_1$ , and thus  $\|\mathbf{p}^*\| \leq \frac{\bar{r}}{\underline{d}}$ .

Similarly, by the optimality of  $\mathbf{p}_t^*$  and its associated  $\mathbf{d}_t \in \mathcal{D}$  in Lemma A.1, we know

$$(\underline{d} - \delta) \|\mathbf{p}_t\|_1 \leq \mathbf{d}_t^\top \mathbf{p}_t \leq \mathbb{E}[r] \leq \bar{r},$$

so we get the bound as  $\|\mathbf{p}_t\| \leq \frac{\bar{r}}{\underline{d} - \delta}$ .

**Proof of Lemma A.4.** According to the latest results in (Jiang et al., 2022), the online dual price  $\mathbf{p}_t$  achieves a sublinear convergence  $\mathcal{O}(\frac{1}{\sqrt{t}})$  to  $\mathbf{p}_t^*$ . Since

$$\begin{aligned} \mathbf{p}^* &\in \arg \min_{\mathbf{p} \geq \mathbf{0}} f(\mathbf{p}) := \mathbf{d}^\top \mathbf{p} + \mathbb{E}[(r - \mathbf{a}^\top \mathbf{p})^+], \\ \mathbf{p}_t^* &\in \arg \min_{\mathbf{p}_t \geq \mathbf{0}} f_{\mathbf{d}_t}(\mathbf{p}_t) := \mathbf{d}_t^\top \mathbf{p}_t + \mathbb{E}[(r - \mathbf{a}^\top \mathbf{p}_t)^+], \end{aligned}$$

By Lemma 12 in (Li & Ye, 2022), we have  $\|\mathbf{p}_t^* - \mathbf{p}^*\|_2^2 \leq \frac{1}{\nu^2 \lambda^2} \|\mathbf{d}_t - \mathbf{d}\|_2^2$ .

### A.3. First-order Analysis

We present the classic first-order method (Li et al., 2020) in Algorithm 4 and the first-order method with restart strategy (Gao et al., 2024) in Algorithm 5. We also show their properties of boundedness and dual convergence results.

---

#### Algorithm 4 First-Order Online algorithm

---

**Input:** total resource  $\mathbf{b}$ , time horizon  $T$ , average resource  $\mathbf{d} = \mathbf{b}/T$ , and initial dual price  $\mathbf{p}_1 = \mathbf{0}$ .

**for**  $t = 1$  to  $T$  **do**

    Observe  $(r_t, \mathbf{a}_t)$  and make decision  $x_t$  based on rule (3).

    Update learning rate  $\alpha_t = \mathcal{O}(1/\sqrt{T})$ .

    Compute subgradient and obtain dual price  $\mathbf{p}_{t+1}$ :

$$\mathbf{p}_{t+1} = \arg \min_{\mathbf{p} \geq \mathbf{0}} (\mathbf{d} - \mathbf{a}_t x_t)^\top \mathbf{p} + \frac{1}{2\alpha_t} \|\mathbf{p} - \mathbf{p}_t\|^2$$

**end**

---

**Algorithm 5** First-Order Restart algorithm

---

**Input:** total resource  $\mathbf{b}$ , time horizon  $T$ , average resource  $\mathbf{d} = \mathbf{b}/T$ , and initial dual price  $\mathbf{p}_1 = \mathbf{0}$ ,  $\mathbf{p}_1^L = \mathbf{0}$

**for**  $t = 1$  to  $T$  **do**

Observe  $(r_t, \mathbf{a}_t)$  and make decision  $x_t$  based on rule (4)

Update learning rate  $\alpha_t = \begin{cases} \mathcal{O}(1/T^{1/3}) & t \leq T^{2/3} \\ \mathcal{O}(1/T^{2/3}) & t > T^{2/3} \end{cases}$

Compute subgradient and obtain dual price  $\mathbf{p}_{t+1}$ :

$$\mathbf{p}_{t+1} = \underset{\mathbf{p} \geq \mathbf{0}}{\operatorname{argmin}} (\mathbf{d} - \mathbf{a}_t x_t)^\top \mathbf{p} + \frac{1}{2\alpha_t} \|\mathbf{p} - \mathbf{p}_t\|^2$$

Run subgradient method with stepsize  $\alpha_t = \mathcal{O}(1/t)$  and update  $\{\mathbf{p}_t^L\}$ . At  $t = T^{2/3}$ , restart  $\mathbf{p}_{T^{2/3}}^L = \mathbf{p}_{T^{2/3}}^L$ .

**end**

---

**Lemma A.5** (Boundedness of first-order result). *The online dual price  $\mathbf{p}_t$  from Algorithm 4 is bounded as*

$$\|\mathbf{p}_t\| \leq \frac{2\bar{r} + m(\bar{a} + \bar{d})^2}{\bar{d}} + m(\bar{a} + \bar{d}).$$

**Lemma A.6** (Dual convergence of First-order algorithm). *Under Assumptions 2.1 and 2.2,  $\mathbf{p}_t$  represents the online solution from first-order method, if  $\alpha_t < \nu\lambda$ , the subgradient updates satisfy the following recursion rule:*

$$\mathbb{E} [\|\mathbf{p}_{t+1} - \mathbf{p}^*\|^2] \leq (1 - \alpha_t \nu \lambda) \|\mathbf{p}_t - \mathbf{p}^*\|^2 + \alpha_t^2 m(\bar{a} + \bar{d})^2. \quad (15)$$

**Case 1.** if  $\alpha_t \equiv \alpha < \frac{1}{\nu\lambda}$ , then there exists a constant  $C_{fo} = \frac{\bar{p}^2 + m(\bar{a} + \bar{d})^2}{\nu\lambda}$  such that

$$\mathbb{E} [\|\mathbf{p}_t - \mathbf{p}^*\|^2] \leq C_{fo} \left( \frac{1}{\alpha t} + \alpha \right). \quad (16)$$

**Case 2.** if  $\alpha_t = \frac{2}{\nu\lambda(t+1)}$ , then there exists a constant  $C_{fo} = \frac{4m(\bar{a} + \bar{d})^2}{\nu^2\lambda^2}$  such that

$$\mathbb{E} [\|\mathbf{p}_t - \mathbf{p}^*\|^2] \leq \frac{C_{fo}}{t}. \quad (17)$$

This lemma shows that the first-order method (Gao et al., 2024) guarantees the convergence of  $\mathbf{p}_t$  to  $\mathbf{p}^*$ . It also highlights a key trade-off in choosing the learning rate  $\alpha_t$ , an aspect that will be central to our later optimality analysis. It also highlights a key trade-off in choosing the learning rate  $\alpha_t$ , an aspect that will be central to our later optimality analysis. With this groundwork in place, we now move on to bounding the total regret of our algorithm.

**Proof of Lemma A.5.** As our initial choice  $\mathbf{p}_1 = \mathbf{0}$ , according to Lemma 1 in (Li et al., 2020), we get the above result. In addition, by Lemma B.1 in (Gao et al., 2024), we have

$$\|\mathbf{p}_t\| \leq \frac{\bar{r}}{\bar{d}} + \frac{m(\bar{a} + \bar{d})^2 \alpha_t}{2\bar{d}} + \alpha_t \sqrt{m}(\bar{a} + \bar{d}).$$

**Proof of Lemma A.6.** Based on the updated rule in Algorithm 4, we derive:

$$\begin{aligned} \|\mathbf{p}_{t+1} - \mathbf{p}^*\|^2 &\leq \|\mathbf{p}_t - \alpha_t(\mathbf{d} - \mathbf{a}_t x_t) - \mathbf{p}^*\|^2 \\ &= \|\mathbf{p}_t - \mathbf{p}^*\|^2 - 2\alpha_t \langle \mathbf{d} - \mathbf{a}_t x_t, \mathbf{p}_t - \mathbf{p}^* \rangle + \alpha_t^2 \|\mathbf{d} - \mathbf{a}_t x_t\|^2 \\ &\leq \|\mathbf{p}_t - \mathbf{p}^*\|^2 - 2\alpha_t \langle \mathbf{d} - \mathbf{a}_t x_t, \mathbf{p}_t - \mathbf{p}^* \rangle + \alpha_t^2 m(\bar{a} + \bar{d})^2. \end{aligned}$$

1) With convexity of  $f$  and  $\mathbb{E}[\mathbf{d} - \mathbf{a}_t x_t] \in \partial f(\mathbf{p}_t)$ , we have:

$$f(\mathbf{p}^*) \geq f(\mathbf{p}_t) + \langle \mathbf{d} - \mathbf{a}_t x_t, \mathbf{p}_t - \mathbf{p}^* \rangle.$$

2) With quadratic regularity of  $f$  as in Lemma A.2, we have:

$$f(\mathbf{p}_t) \geq f(\mathbf{p}^*) + \nabla f(\mathbf{p}^*)^\top (\mathbf{p}_t - \mathbf{p}^*) + \frac{\nu\lambda}{2} \|\mathbf{p}_t - \mathbf{p}^*\|^2,$$

which indicates  $f(\mathbf{p}_t) - f(\mathbf{p}^*) \geq \frac{\nu\lambda}{2} \|\mathbf{p}_t - \mathbf{p}^*\|^2$ .

Combine the above results and take expectation conditioned on history information  $\{(r_j, \mathbf{a}_j), j \leq t\}$ , we obtain:

$$\begin{aligned} \mathbb{E}\|\mathbf{p}_{t+1} - \mathbf{p}^*\|^2 &\leq \|\mathbf{p}_t - \mathbf{p}^*\|^2 - 2\alpha_t \mathbb{E}[\langle \mathbf{d} - \mathbf{a}_t x_t, \mathbf{p}_t - \mathbf{p}^* \rangle] + \alpha_t^2 m(\bar{a} + \bar{d})^2 \\ &\leq \|\mathbf{p}_t - \mathbf{p}^*\|^2 - 2\alpha_t (f(\mathbf{p}_t) - f(\mathbf{p}^*)) + \alpha_t^2 m(\bar{a} + \bar{d})^2 \\ &\leq \|\mathbf{p}_t - \mathbf{p}^*\|^2 - \alpha_t \nu \lambda \|\mathbf{p}_t - \mathbf{p}^*\|^2 + \alpha_t^2 m(\bar{a} + \bar{d})^2 \\ &= (1 - \alpha_t \nu \lambda) \|\mathbf{p}_t - \mathbf{p}^*\|^2 + \alpha_t^2 m(\bar{a} + \bar{d})^2. \end{aligned}$$

This proves (15) with the general case of learning rate  $\alpha_t$ .

**Case 1.** When  $\alpha_t = \alpha < \frac{1}{\nu\lambda}$  is a constant, we take recursion of (15). Note that  $(1 - \nu\lambda\alpha)^t < 1/\nu\lambda\alpha t$ , we get:

$$\begin{aligned} \mathbb{E}\|\mathbf{p}_{t+1} - \mathbf{p}^*\|^2 &\leq (1 - \alpha\nu\lambda)^t \|\mathbf{p}_1 - \mathbf{p}^*\|^2 + \sum_{j=0}^{t-1} \alpha^2 m(\bar{a} + \bar{d})^2 (1 - \alpha\nu\lambda)^j \\ &\leq \frac{\|\mathbf{p}_1 - \mathbf{p}^*\|^2}{\nu\lambda\alpha t} + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} \alpha. \end{aligned}$$

Since all feasible  $\mathbf{p} \in \mathcal{V}_{\mathbf{p}}$  is bounded, let  $C_{fo} = \frac{\bar{p}^2 + m(\bar{a} + \bar{d})^2}{\nu\lambda}$ , we are able to obtain (16).

**Case 2.** When  $\alpha_t = \frac{2}{\nu\lambda(t+1)}$ , for any  $j \leq t$ , Lemma A.6 gives us

$$\mathbb{E}\|\mathbf{p}_{j+1} - \mathbf{p}^*\|^2 \leq \frac{j-1}{j+1} \|\mathbf{p}_j - \mathbf{p}^*\|^2 + \frac{4m(\bar{a} + \bar{d})^2}{\nu^2 \lambda^2 (j+1)^2}.$$

Re-arranging the above inequality, we have

$$\begin{aligned} (j+1)^2 \mathbb{E}\|\mathbf{p}_{j+1} - \mathbf{p}^*\|^2 &\leq (j^2 - 1) \|\mathbf{p}_j - \mathbf{p}^*\|^2 + \frac{4m(\bar{a} + \bar{d})^2}{\nu^2 \lambda^2}, \\ \text{then } (j+1)^2 \mathbb{E}\|\mathbf{p}_{j+1} - \mathbf{p}^*\|^2 - j^2 \|\mathbf{p}_j - \mathbf{p}^*\|^2 &\leq \frac{4m(\bar{a} + \bar{d})^2}{\nu^2 \lambda^2}. \end{aligned}$$

Then by telescoping from  $j = 2$  to  $t$ , we have

$$\sum_{j=1}^t (j+1)^2 \mathbb{E}\|\mathbf{p}_{j+1} - \mathbf{p}^*\|^2 - j^2 \|\mathbf{p}_j - \mathbf{p}^*\|^2 = (t+1)^2 \mathbb{E}\|\mathbf{p}_{t+1} - \mathbf{p}^*\|^2 - \|\mathbf{p}_1 - \mathbf{p}^*\|^2 \leq \frac{4m(\bar{a} + \bar{d})^2 t}{\nu^2 \lambda^2}$$

which then gives us

$$\begin{aligned} \mathbb{E}\|\mathbf{p}_{t+1} - \mathbf{p}^*\|^2 &\leq \frac{\|\mathbf{p}_1 - \mathbf{p}^*\|^2}{(t+1)^2} + \frac{4m(\bar{a} + \bar{d})^2 t}{\nu^2 \lambda^2 (t+1)^2}, \\ \text{thus } \mathbb{E}\|\mathbf{p}_t - \mathbf{p}^*\|^2 &\leq \frac{\|\mathbf{p}_1 - \mathbf{p}^*\|^2}{t^2} + \frac{4m(\bar{a} + \bar{d})^2}{\nu^2 \lambda^2 t}. \end{aligned}$$

As  $\mathbf{p}_1 = \mathbf{0}$  and  $\mathbf{p}^*$  is bounded, taking  $C_{fo} = \frac{4m(\bar{a} + \bar{d})^2}{\nu^2 \lambda^2}$  completes the proof for (17).

## B. Algorithm Design and Analysis

In this section, we present a cohesive mathematical framework to integrate the LP-based method and first-order method by developing a new performance metric. We also provide some unique properties for our algorithm design of parallel multi-phase structure.

As stated in Algorithm 1, we construct independent paths for online learning and online decision-making conducted simultaneously. Specifically, for the online learning path, we use the LP-based method to resolve updated linear programs at a fixed frequency and send this result to the decision-making path; For the online decision-making path, we apply the first-order method in the initial and final batches, and use the latest dual price for decision-making in the intermediate resolving intervals. This process is illustrated in Figure 1.

### B.1. Regret Analysis

First, we construct an upper bound for the offline optimal objective value. The challenge comes from the intractable dependency of constraints on objective value under the online setting. We tackle this issue by introducing a Lagrangian function to integrate constraints into the objective and balance revenue maximization with constraint satisfaction. The formulation is stated as follows.

**Lemma B.1** (Lagrangian Upper Bound). *Under Assumptions (2.1) and (2.2), define the deterministic Lagrangian dual function as*

$$\ell(\mathbf{p}) := \mathbb{E} [rI(r > \mathbf{a}^\top \mathbf{p}) + (\mathbf{d} - \mathbf{a}I(r > \mathbf{a}^\top \mathbf{p}))^\top \mathbf{p}^*].$$

Then for any feasible  $\mathbf{p} \in \mathcal{V}_{\mathbf{p}}$ , we have:

$$\begin{aligned} (a) \quad & \mathbb{E} \left[ \sum_{t=1}^T r_t x_t^* \right] \leq T \ell(\mathbf{p}^*), \\ (b) \quad & \ell(\mathbf{p}^*) - \ell(\mathbf{p}) \leq \mu \bar{a}^2 \|\mathbf{p} - \mathbf{p}^*\|^2. \end{aligned} \tag{18}$$

**Lemma B.2** (Dual Price Boundedness). *Under Assumption 2.1, the online and offline optimal dual prices are bounded respectively by  $\|\mathbf{p}^*\| \leq \frac{\bar{r}}{\underline{d}}$  and*

$$\|\mathbf{p}_t\| \leq \bar{p} := \max \left( \frac{\bar{r}}{\underline{d} - \delta}, \frac{2\bar{r} + m(\bar{a} + \bar{d})^2}{\underline{d}} + m(\bar{a} + \bar{d}) \right).$$

This lemma establishes that the optimal dual prices remain bounded. Our algorithm maintains these bounds because if  $\mathbf{p}_t$  grows large, the algorithm responds by accepting more orders, which in turn reduces  $\mathbf{p}_{t+1}$ . This self-correcting mechanism keeps the dual prices within these limits. With this groundwork in place, we now move on to derive the upper bound for the total performance of Algorithm 1.

**Theorem B.3** (Decomposition of Regret). *Under Assumptions (2.1) and (2.2), let  $C_r = \max\{\frac{\bar{r}}{\underline{d}}, \frac{\nu\lambda\bar{a}^2}{2}\}$ , we derive an upper bound for the regret  $r(\mathbf{x})$  as:*

$$r(\mathbf{x}) \leq C_r \left[ \mathbb{E} \left\| \left( \mathbf{b} - \sum_{t=1}^T \mathbf{a}_t x_t \right)^{B^+} \right\| + \sum_{t=1}^T \mathbb{E} \|\mathbf{p}_t - \mathbf{p}^*\|^2 \right]. \tag{19}$$

**Proof of Lemma B.1.** This lemma is proved as Lemma 3 in (Li & Ye, 2022) by strong duality and optimality of  $\mathbf{p}_T^*$ .

**Proof of Lemma B.2.** By the boundedness results in Lemma A.3 and Lemma A.5, define

$$\bar{p} := \max \left( \frac{\bar{r}}{\underline{d} - \delta}, \frac{2\bar{r} + m(\bar{a} + \bar{d})^2}{\underline{d}} + m(\bar{a} + \bar{d}) \right),$$

then we complete the proof.

**Proof of Theorem B.3.** By the definition of  $\ell(\mathbf{p})$ , we derive the online objective as

$$\begin{aligned}\mathbb{E}\left[\sum_{t=1}^T r_t x_t\right] &= \sum_{t=1}^T \mathbb{E}[\mathbb{E}[r_t x_t | \mathbf{p}_t]] \\ &= \sum_{t=1}^T \mathbb{E}[\ell(\mathbf{p}_t) - (\mathbf{d} - \mathbf{a}_t x_t)^\top \mathbf{p}^*].\end{aligned}$$

Then by Lemma B.1, we derive the upper bound for the regret  $r(\mathbf{x})$  as defined in (6) to be:

$$\begin{aligned}r(\mathbf{x}) &= \mathbb{E}\left[\sum_{t=1}^T r_t x_t^* - r_t x_t\right] \\ &\leq T\ell(\mathbf{p}^*) - \sum_{t=1}^T \mathbb{E}[\ell(\mathbf{p}_t) - (\mathbf{d} - \mathbf{a}_t x_t)^\top \mathbf{p}^*] \\ &= \sum_{t=1}^T \mathbb{E}[(\mathbf{d} - \mathbf{a}_t x_t)^\top \mathbf{p}^*] + \sum_{t=1}^T \mathbb{E}[\ell(\mathbf{p}^*) - \ell(\mathbf{p}_t)] \\ &\leq \mathbb{E}\left[\left(\mathbf{b} - \sum_{t=1}^T \mathbf{a}_t x_t\right)^\top \mathbf{p}^*\right] + \frac{\mu \bar{a}^2}{2} \sum_{t=1}^T \mathbb{E}\|\mathbf{p}_t - \mathbf{p}^*\|^2 \\ &\leq \|\mathbf{p}^*\| \cdot \mathbb{E}\left\|\left(\mathbf{b} - \sum_{t=1}^T \mathbf{a}_t x_t\right)^{B^+}\right\| + \frac{\mu \bar{a}^2}{2} \sum_{t=1}^T \mathbb{E}\|\mathbf{p}_t - \mathbf{p}^*\|^2\end{aligned}$$

where  $(\cdot)^{B^+}$  denotes the positive part only for binding constraints.

As  $\|\mathbf{p}^*\| \leq \frac{\bar{r}}{\bar{d}}$  is bounded, taking the constant  $C_r = \max\{\frac{\bar{r}}{\bar{d}}, \frac{\mu \bar{a}^2}{2}\}$  completes the proof.

## B.2. Total Performance Analysis

We consider the constraint violation as  $v(\mathbf{x}) = \|(\mathbf{A}\mathbf{x} - \mathbf{b})^+\|$  in (7). Then the total performance is the combination of regret and constraint violation. Since Theorem B.3 holds for any unknown distribution  $\mathcal{P} \in \Xi$ , we have:

$$\begin{aligned}\Delta_T &= \sup_{\mathcal{P} \in \Xi} \mathbb{E}_{\mathcal{P}}[r(\mathbf{x}) + v(\mathbf{x})] \\ &\leq C_r \left[ \mathbb{E}\left\|\left(\mathbf{b} - \sum_{t=1}^T \mathbf{a}_t x_t\right)^{B^+}\right\| + \sum_{t=1}^T \mathbb{E}\|\mathbf{p}_t - \mathbf{p}^*\|^2 \right] + \mathbb{E}\left\|\left(\sum_{t=1}^T \mathbf{a}_t x_t - \mathbf{b}\right)^+\right\| \\ &\leq C_r \left[ \sum_{t=1}^T \mathbb{E}\|\mathbf{p}_t - \mathbf{p}^*\|^2 + \mathbb{E}\|(\mathbf{b} - \mathbf{A}\mathbf{x})^{B^+}\| + \mathbb{E}\|(\mathbf{A}\mathbf{x} - \mathbf{b})^+\| \right].\end{aligned}\tag{20}$$

This new framework is adaptive to all OLP methods. We proceed with this structure to analyze various OLP methods.

**Theorem B.4** (Horizon Division). *Based on the change of methods, we separate the horizon into three intervals  $T_1, T_2$ , and  $T_3$ , and reorganize the total performance for analysis to be*

$$\Delta_T \leq C_r [\Delta_{T_1} + \Delta_{T_2} + \Delta_{T_3}].\tag{21}$$

where  $\Delta_{T_1}, \Delta_{T_2}$ , and  $\Delta_{T_3}$  are the performances for each interval.

**Proof of Theorem B.4.** We split the total horizon into three intervals of initial batch  $T_1 = [0, f]$ , intermediate process  $T_2 = [f, kf]$ , and final batch  $T_3 = [kf, T]$ . We use the first-order method for the first and final batches and the LP-based method for intermediate processes. Notice that for the initial resolving batch, we maintain the classical analysis for the



first-order method since it has not reached the mixture with the LP-based method. Then based on (20), we define

$$\begin{aligned}
 \Delta_{T_1} &:= \sum_{t=1}^f \mathbb{E}[r_t x_t^* - r_t x_t] + \mathbb{E} \left\| \left( \sum_{t=1}^f \mathbf{a}_t x_t - f \mathbf{d} \right)^+ \right\| \\
 \Delta_{T_2} &:= \sum_{t=f}^{kf} \mathbb{E} \|\mathbf{p}_t - \mathbf{p}^*\|^2 + \mathbb{E} \left\| \left( (k-1)f \mathbf{d} - \sum_{t=f}^{kf} \mathbf{a}_t x_t \right)^{B^+} \right\| + \mathbb{E} \left\| \left( \sum_{t=f}^{kf} \mathbf{a}_t x_t - (k-1)f \mathbf{d} \right)^+ \right\| \\
 \Delta_{T_3} &:= \sum_{t=kf}^T \mathbb{E} \|\mathbf{p}_t - \mathbf{p}^*\|^2 + \mathbb{E} \left\| \left( (T-kf) \mathbf{d} - \sum_{t=kf}^T \mathbf{a}_t x_t \right)^{B^+} \right\| + \mathbb{E} \left\| \left( \sum_{t=kf}^T \mathbf{a}_t x_t - (T-kf) \mathbf{d} \right)^+ \right\|
 \end{aligned} \tag{22}$$

Then we have

$$\Delta_T \leq C_r [\Delta_{T_1} + \Delta_{T_2} + \Delta_{T_3}]$$

since  $\|(a+b)^+\| \leq \|a^+\| + \|b^+\|$  for any  $a, b$ .

### B.3. Regret for LP-based Method

In this section, we are going to bound  $\Delta_{T_2}$  for the regret from LP-based method. To provide a tighter analysis, we examine the real-time average resource capacity  $\mathbf{d}_t$  instead of the original process  $\mathbf{b}_t$ . Specifically, by Lemma A.1, we constraint  $\mathbf{d}_t$  within a feasible range  $\mathcal{D}$ . If  $\mathbf{d}_t \notin \mathcal{D}$ , we set the dual price to zero to accept all subsequent orders. We analyze the remaining resources and constraint violation due to this strategy. Note that this is required solely for rigor analysis purposes and is not necessary when running the algorithm.

**Lemma B.5** (Dynamics of Resource Usage). *Under Assumption 2.1 and 2.2, there exists a constant  $C > 0$  depending on  $\bar{d}, \bar{a}, m, \nu, \lambda_1, \mu$ , and  $C_{lp}$  such that*

$$\sum_{t=f}^{kf} \mathbb{E}[(d_{i,t} - d_i)^2] \leq C \log(k). \tag{23}$$

**Theorem B.6** (Regret of Intermediate Process). *Under Assumption 2.1 and 2.2, following the result in (22), we prove the regret satisfies*

$$\Delta_{T_2} \leq \log\left(\frac{T}{f}\right) = \mathcal{O}(\log(k)). \tag{24}$$

**Proof of Lemma B.5.** During the re-solving process, we denote each interval to be  $[jf, (j+1)f]$  where  $j = 1, 2, \dots, k$ . As the resource usage follows  $\mathbf{b}_{t+1} = \mathbf{b}_t - \mathbf{a}_{t+1} I(r_{t+1} > \mathbf{a}_{t+1}^\top \mathbf{p}_{t+1})$ , normalizing both sides, we derive the update of average resource consumption to be:

$$d_{i,(j+1)f} = d_{i,jf} + \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f}.$$

Subtracting  $d_i$  on both sides, it becomes

$$\begin{aligned}
 d_{i,(j+1)f} - d_i &= d_{i,jf} - d_i + \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f} \\
 &= d_{i,jf} - d_i + \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*)}{T - (j+1)f} \\
 &\quad + \frac{\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f}.
 \end{aligned}$$

Taking expectations of squares, we have

$$\begin{aligned}
 \mathbb{E}(d_{i,(j+1)f} - d_i)^2 &= \mathbb{E}(d_{i,jf} - d_i)^2 + \mathbb{E} \left[ \frac{(\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*))^2}{(T - (j+1)f)^2} \right] \\
 &\quad + \mathbb{E} \left[ \frac{(\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}))^2}{(T - (j+1)f)^2} \right] \\
 &\quad + 2\mathbb{E} \left[ (d_{i,jf} - d_i) \left( \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*)}{T - (j+1)f} \right) \right] \\
 &\quad + 2\mathbb{E} \left[ (d_{i,jf} - d_i) \left( \frac{\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f} \right) \right] \\
 &\quad + 2\mathbb{E} \left[ \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*)}{T - (j+1)f} \cdot \frac{\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f} \right].
 \end{aligned} \tag{25}$$

By Lemma D.1, we obtain the recursion relation as

$$\mathbb{E}(d_{i,(j+1)f} - d_i)^2 \leq \mathbb{E}(d_{i,jf} - d_i)^2 + \frac{C_{rec}}{(k-j-1)^2 f} + \frac{4\mu\bar{a}^2\sqrt{C_{lp}}}{(k-j-1)\sqrt{(j+1)f}} \sqrt{\mathbb{E}[(d_{i,jf} - d_i)^2]} \tag{26}$$

where  $C_{rec} > 0$  is a constant defined in Lemma D.1.

Then according to Lemma D.2, take  $C = 12 \max\{C_{rec}, 16\mu^2\bar{a}^4 C_{lp}\}$ , we solve recursion (26) and obtain the upper bound of total deviation from original  $d_i$  in the re-solving process to be:

$$\sum_{j=1}^k \mathbb{E}[(d_{i,jf} - d_i)^2] \leq \frac{C}{f} \log(k). \tag{27}$$

Therefore, we sum the whole re-solving process and obtain:

$$\begin{aligned}
 \sum_{t=f}^{kf} \mathbb{E}[(d_{i,t} - d_i)^2] &= \sum_{j=1}^k \sum_{\ell=jf}^{(j+1)f} \mathbb{E}[(d_{i,\ell} - d_i)^2] \\
 &= f \cdot \frac{C}{f} \log(k) \quad (\text{by result in (27)}) \\
 &\leq C \log(k).
 \end{aligned} \tag{28}$$

This completes the proof.

**Proof of Theorem B.6.** We analyze the three terms in (22) respectively.

1. For dual convergence, by Lemma A.4, we have

$$\begin{aligned}
 \sum_{t=f}^{kf} \mathbb{E}\|\mathbf{p}_t - \mathbf{p}^*\|^2 &= \mathbb{E} \left[ \sum_{j=1}^{k-1} \sum_{t=jf+1}^{(j+1)f} \|\mathbf{p}_t - \mathbf{p}^*\|_2^2 \right] \\
 &= \sum_{j=1}^{k-1} f \mathbb{E} \left[ \|\mathbf{p}_{(j+1)f} - \mathbf{p}^*\|_2^2 \right] \\
 &\leq \sum_{j=1}^{k-1} f \mathbb{E} \left[ \|\mathbf{p}_{(j+1)f} - \mathbf{p}_{jf}^*\|_2^2 + \|\mathbf{p}_{jf}^* - \mathbf{p}^*\|_2^2 \right]
 \end{aligned}$$

$$\begin{aligned}
 &\leq \sum_{j=1}^{k-1} f \left[ \frac{C_{lp}}{jf} + \frac{1}{\nu^2 \lambda^2} \mathbb{E}[(d_{jf} - d)^2] \right] \quad (\text{by Lemma A.4}) \\
 &\leq C_{lp} \log(k) + \frac{mC}{\nu^2 \lambda^2} \log(k). \quad (\text{by Lemma B.5}).
 \end{aligned} \tag{29}$$

2. For the constraint violation, by Chebyshev's inequality, we bound the probability as

$$\begin{aligned}
 \sum_{t=f}^{kf} \mathbb{P}(|d_{i,t} - d_i| \leq \delta) &\leq \sum_{t=f}^{kf} \frac{\mathbb{E}[(d_{i,t} - d_i)^2]}{\delta^2} \\
 &\leq \frac{C}{\delta^2} \log(k). \quad (\text{by Lemma B.5})
 \end{aligned} \tag{30}$$

Since our strategy is to automatically set  $p_t = 0$  once  $d_t \notin \mathcal{D}$ , we will accept all the subsequent orders according to our decision condition (3). Denote  $R$  as the total number of orders processed specifically due to this rule, then we have:

$$\mathbb{E}[R] \leq \sum_{t=f}^{kf} \mathbb{P}(|d_{i,t} - d_i| \leq \delta) \leq \frac{C}{\delta^2} \log(k) \tag{31}$$

which indicates that the resource violation is at most

$$\mathbb{E} \left\| \left( \sum_{t=f}^{kf} \mathbf{a}_t x_t - (k-1)f\mathbf{d} \right)^+ \right\| \leq \mathbb{E}[\bar{a}R] \leq \frac{C\bar{a}}{\delta^2} \log(k). \tag{32}$$

3. For the remaining resource, we derive its upper bound by considering the opposite extreme case of our strategy. If we reject all those  $R$  orders, then no resource is used on them. Thus, by (31), the positive projection of remaining resource at  $t = kf$  is at most:

$$\mathbb{E} \left\| \left( (k-1)f\mathbf{d} - \sum_{t=f}^{kf} \mathbf{a}_t x_t \right)^{B+} \right\| \leq \mathbb{E}[(d_i + \delta)R] \leq \frac{C(\bar{d} + \delta)}{\delta^2} \log(k). \tag{33}$$

Combining the results of (29), (32), and (33), we derive the final result as

$$\Delta_{T_2} \leq \left( C_{lp} + \frac{mC}{\nu^2 \lambda^2} + (\bar{a} + \bar{d} + \delta) \frac{C}{\delta^2} \right) \log(k). \tag{34}$$

This completes the proof.

#### B.4. Regret for First-Order Method

In this section, we are going to bound  $\Delta_{T_1}$  and  $\Delta_{T_3}$  for the regret from the first-order method. The key point is to analyze the connection between dual prices based on the gradient-updated rule. Here  $p_t = p_t^{fo}$ .

**Theorem B.7** (Regret of Initial Batch). *Under Assumption 2.1 and 2.2, we have*

$$\Delta_{T_1} \leq \left( \frac{m(\bar{a} + \bar{d})^2}{2} + \frac{2\bar{r} + m(\bar{a} + \bar{d})^2}{\underline{d}} + m(\bar{a} + \bar{d}) \right) \sqrt{f}. \tag{35}$$

**Lemma B.8** (Constraint Violation). *Under Assumption 2.1 and 2.2, take the learning rate  $\alpha_t = \alpha$ , we bound the constraint violation for the last batch as*

$$\mathbb{E} \left\| \left( \sum_{t=kf}^T \mathbf{a}_t x_t - (T - kf)\mathbf{d} \right)^+ \right\| \leq \left( \frac{C_{lp}}{\sqrt{kf}} + \frac{C_b}{T^2} \right) \cdot \frac{1}{\alpha} + \frac{C_b}{\sqrt{(T - kf)kf}} \cdot \frac{1}{\alpha\sqrt{\alpha}} + \frac{C_b}{\sqrt{\alpha}}. \tag{36}$$

**Theorem B.9** (Regret of Final Batch). *Under Assumption 2.1 and 2.2, we have*

$$\Delta_{T_3} \leq \left( \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} + \frac{2C_{lp}}{\sqrt{k}} + \frac{C_{lp} \log(f)}{\nu\lambda k} + \frac{4C_b}{\sqrt{k}} + 2C_b \right) \cdot f^{1/3}. \quad (37)$$

**Lemma B.10** (Warm-Start First-Order Regret). *Assume that we have a warm start for the initial dual price  $\mathbf{p}_0$  in the first batch, which satisfies  $\|\mathbf{p}_0 - \mathbf{p}^*\| \leq f^{-1/3}$ . Then the regret of the first batch follows*

$$\Delta_{T_1 \text{warm}} \leq \left( \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} + 4C_w \right) \cdot f^{1/3} + \log(f). \quad (38)$$

**Proof of Theorem B.7.** We analyze  $\Delta_{T_1}$  in (22) by parts. By Lemma B.3, B.4 in (Gao et al., 2024), we have

$$\begin{aligned} \sum_{t=1}^f \mathbb{E}[r_t x_t^* - r_t x_t] &\leq \frac{m(\bar{a} + \bar{d})^2}{2} \alpha f, \\ \mathbb{E} \left\| \left( \sum_{t=1}^f \mathbf{a}_t x_t - f \mathbf{d} \right)^+ \right\| &\leq \frac{1}{\alpha} \left[ \frac{2\bar{r} + m(\bar{a} + \bar{d})^2}{\underline{d}} + m(\bar{a} + \bar{d}) \right]. \end{aligned}$$

To achieve the tight upper bound, we select the optimal step size  $\alpha = \frac{1}{\sqrt{f}}$ , which then gives us the desired result.

**Proof of Lemma B.8.** With the update rule of first-order method in Algorithm 4, we have

$$\begin{aligned} \mathbf{p}_{t+1} &= [\mathbf{p}_t - \alpha(\mathbf{d} - \mathbf{a}_t x_t)]^+ \geq \mathbf{p}_t - \alpha(\mathbf{d} - \mathbf{a}_t x_t), \\ \text{which gives us } \mathbf{a}_t x_t - \mathbf{d} &\leq \frac{1}{\alpha}(\mathbf{p}_{t+1} - \mathbf{p}_t). \end{aligned}$$

Summarizing on both sides and applying the telescoping, we derive

$$\sum_{t=kf}^T (\mathbf{a}_t x_t - \mathbf{d}) \leq \frac{1}{\alpha} \sum_{t=kf}^T (\mathbf{p}_{t+1} - \mathbf{p}_t) = \frac{1}{\alpha} (\mathbf{p}_{T+1} - \mathbf{p}_{kf}).$$

By Lemma D.3, take  $C_b = \max\{\frac{C_{lp}}{\sqrt{\nu\lambda}}, \frac{m(\bar{a} + \bar{d})}{\sqrt{\nu\lambda}}, \bar{p}\}$ , we have

$$\mathbb{E} \|\mathbf{p}_{T+1} - \mathbf{p}^*\| \leq C_b \left[ \frac{1}{\sqrt{(T - kf)kf}} \cdot \frac{1}{\sqrt{\alpha}} + \sqrt{\alpha} + \frac{1}{T^2} \right]. \quad (39)$$

Thus, according to Lemma A.4 and (39), the expectation of constraint violation satisfies

$$\begin{aligned} \mathbb{E} \left\| \left( \sum_{t=kf}^T (\mathbf{a}_t x_t - \mathbf{d}) \right)^+ \right\| &\leq \mathbb{E} \left\| \sum_{t=kf}^T (\mathbf{a}_t x_t - \mathbf{d}) \right\| \\ &\leq \frac{1}{\alpha} \mathbb{E} \|\mathbf{p}_{T+1} - \mathbf{p}_{kf}\| \\ &\leq \frac{1}{\alpha} \mathbb{E} [\|\mathbf{p}_{kf} - \mathbf{p}^*\| + \|\mathbf{p}_{T+1} - \mathbf{p}^*\|] \\ &\leq \frac{1}{\alpha} \left( \frac{C_{lp}}{\sqrt{kf}} + C_b \left[ \frac{1}{\sqrt{(T - kf)kf}} \cdot \frac{1}{\sqrt{\alpha}} + \sqrt{\alpha} + \frac{1}{T^2} \right] \right) \\ &\leq \frac{C_{lp}}{\sqrt{kf}} \cdot \frac{1}{\alpha} + \frac{C_b}{\sqrt{(T - kf)kf}} \cdot \frac{1}{\alpha\sqrt{\alpha}} + C_b \cdot \frac{1}{\sqrt{\alpha}} + \frac{C_b}{T^2} \cdot \frac{1}{\alpha}. \end{aligned} \quad (40)$$

This completes the proof.

**Proof of Theorem B.9.** We decompose  $\Delta_{T_3}$  into three parts according to (22) and analyze each term respectively. Since the first-order method re-starts from  $\mathbf{p}_{kf}$ , by Lemma A.4, we have:

$$\begin{aligned}
 \sum_{t=kf}^T \mathbb{E} \|\mathbf{p}_t - \mathbf{p}^*\|^2 &\leq \sum_{t=kf}^T \left[ \frac{\|\mathbf{p}_{kf} - \mathbf{p}^*\|^2}{\nu\lambda\alpha t} + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} \alpha \right] \\
 &\leq \frac{\|\mathbf{p}_{kf} - \mathbf{p}^*\|^2}{\nu\lambda\alpha} \log(T - kf) + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} (T - kf) \alpha \\
 &\leq \frac{\|\mathbf{p}_{kf} - \mathbf{p}^*\|^2}{\nu\lambda\alpha} \log(f) + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f \alpha \\
 &\leq \frac{C_{lp} \log(f)}{\nu\lambda kf} \cdot \frac{1}{\alpha} + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f \alpha.
 \end{aligned} \tag{41}$$

Since  $t \in [kf, T]$ , by Proposition 3.3 in (Gao et al., 2024), the dual price is far from the origin and thus the positive projection still equals itself. This indicates that:

$$\begin{aligned}
 \mathbf{p}_{t+1} &= [\mathbf{p}_t - \alpha(\mathbf{d} - \mathbf{a}_t x_t)]^+ = \mathbf{p}_t - \alpha(\mathbf{d} - \mathbf{a}_t x_t), \\
 \text{which then gives us } \mathbf{d} - \mathbf{a}_t x_t &= \frac{1}{\alpha}(\mathbf{p}_t - \mathbf{p}_{t+1}).
 \end{aligned}$$

Summarizing on both sides and applying the telescoping, we derive

$$\sum_{t=kf}^T (\mathbf{d} - \mathbf{a}_t x_t) \leq \frac{1}{\alpha} \sum_{t=kf}^T (\mathbf{p}_t - \mathbf{p}_{t+1}) = \frac{1}{\alpha} (\mathbf{p}_{kf} - \mathbf{p}_{T+1}).$$

Then for the positive binding terms, by Lemma A.4 and (39), we have

$$\begin{aligned}
 \mathbb{E} \left\| \left( \sum_{t=kf}^T (\mathbf{d} - \mathbf{a}_t x_t) \right)^{B^+} \right\| &\leq \mathbb{E} \left\| \sum_{t=kf}^T (\mathbf{d} - \mathbf{a}_t x_t) \right\| \\
 &\leq \frac{1}{\alpha} \mathbb{E} \|\mathbf{p}_{kf} - \mathbf{p}_{T+1}\| \\
 &\leq \frac{1}{\alpha} \mathbb{E} [\|\mathbf{p}_{kf} - \mathbf{p}^*\| + \|\mathbf{p}_{T+1} - \mathbf{p}^*\|] \\
 &\leq \frac{C_{lp}}{\sqrt{kf}} \cdot \frac{1}{\alpha} + \frac{C_b}{\sqrt{(T - kf)kf}} \cdot \frac{1}{\alpha\sqrt{\alpha}} + C_b \cdot \frac{1}{\sqrt{\alpha}} + \frac{C_b}{T^2} \cdot \frac{1}{\alpha}.
 \end{aligned} \tag{42}$$

Combining the results of (41), Lemma B.8, and (42) together, we obtain the final result as:

$$\Delta_{T_3} \leq \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f \alpha + \left( \frac{2C_{lp}}{\sqrt{kf}} + \frac{2C_b}{T^2} + \frac{C_{lp} \log(f)}{\nu\lambda kf} \right) \cdot \frac{1}{\alpha} + \frac{2C_b}{\sqrt{(T - kf)kf}} \cdot \frac{2}{\alpha\sqrt{\alpha}} + \frac{2C_b}{\sqrt{\alpha}}. \tag{43}$$

We select the optimal learning rate  $\alpha = f^{-2/3}$  to minimize  $\Delta_{T_3}$  in (43). We prove this result in cases.

1. If  $T - kf \leq f^{1/3}$ , the regret must be smaller than the length of the batch, so  $\Delta_{T_3} \leq T - kf \leq f^{1/3}$ .
2. If  $T - kf > f^{1/3}$ , then we use this property to bound the term  $\frac{2C_b}{\sqrt{(T - kf)kf}} \cdot \frac{2}{\alpha\sqrt{\alpha}}$ . We derive

$$\begin{aligned}
 \Delta_{T_3} &\leq \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f \cdot f^{-2/3} + \left( \frac{2C_{lp}}{\sqrt{kf}} + \frac{2C_b}{T^2} + \frac{C_{lp} \log(f)}{\nu\lambda kf} \right) \cdot f^{2/3} + \frac{4C_b}{\sqrt{f^{1/3} \cdot kf}} \cdot f + 2C_b \cdot f^{1/3} \\
 &= \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f^{1/3} + \frac{2C_{lp}}{\sqrt{k}} f^{1/6} + \frac{2C_b}{T} + \frac{C_{lp} \log(f)}{\nu\lambda k} f^{1/3} + \frac{4C_b}{\sqrt{k}} \cdot f^{1/3} + 2C_b \cdot f^{1/3} \\
 &\leq \left( \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} + \frac{2C_{lp}}{\sqrt{k}} + \frac{C_{lp} \log(f)}{\nu\lambda k} + \frac{4C_b}{\sqrt{k}} + 2C_b \right) \cdot f^{1/3}.
 \end{aligned} \tag{44}$$

This completes the proof.



**Proof of Lemma B.10.** According to Lemma D.3, with  $\|\mathbf{p}_0 - \mathbf{p}^*\| \leq f^{-1/3}$  and  $\alpha < \frac{1}{\nu\lambda}$ , we have

$$\begin{aligned}\mathbb{E}[\|\mathbf{p}_{f+1} - \mathbf{p}^*\|^2 \mid \mathbf{p}_0] &\leq \frac{\|\mathbf{p}_0 - \mathbf{p}^*\|^2}{\nu\lambda\alpha f} + \frac{\alpha m(\bar{a} + \bar{d})^2}{\nu\lambda}, \\ \mathbb{E}\|\mathbf{p}_{f+1} - \mathbf{p}^*\| &\leq \frac{1}{f^{1/3}\sqrt{\nu\lambda f}} \cdot \frac{1}{\sqrt{\alpha}} + \frac{m(\bar{a} + \bar{d})}{\sqrt{\nu\lambda}} \cdot \sqrt{\alpha} + \frac{\bar{p}}{T^2}.\end{aligned}\quad (45)$$

Then according to Lemma B.8, take  $C_w = \max\{\frac{1}{\nu\lambda}, \frac{m(\bar{a} + \bar{d})}{\sqrt{\nu\lambda}}, \bar{p}\}$ , we derive the constraint violation follows

$$\begin{aligned}\mathbb{E}\left\|\left(\sum_{t=0}^f (\mathbf{a}_t x_t - \mathbf{d})\right)^+\right\| &\leq \mathbb{E}\left\|\sum_{t=0}^f (\mathbf{a}_t x_t - \mathbf{d})\right\| \\ &\leq \frac{1}{\alpha} \mathbb{E}\|\mathbf{p}_{f+1} - \mathbf{p}_0\| \\ &\leq \frac{1}{\alpha} \mathbb{E}[\|\mathbf{p}_0 - \mathbf{p}^*\| + \|\mathbf{p}_{f+1} - \mathbf{p}^*\|] \\ &\leq \frac{1}{\alpha} \left( \frac{1}{f^{1/3}} + C_w \left[ \frac{1}{f^{1/3}\sqrt{f}} \cdot \frac{1}{\sqrt{\alpha}} + \sqrt{\alpha} + \frac{1}{T^2} \right] \right) \\ &\leq \frac{1}{f^{1/3}} \cdot \frac{1}{\alpha} + \frac{C_w}{f^{5/6}} \cdot \frac{1}{\alpha\sqrt{\alpha}} + C_w \cdot \frac{1}{\sqrt{\alpha}} + \frac{C_w}{T^2} \cdot \frac{1}{\alpha}.\end{aligned}\quad (46)$$

Similar to (41) in Theorem B.9, we derive the dual distance as

$$\begin{aligned}\sum_{t=0}^f \mathbb{E}\|\mathbf{p}_t - \mathbf{p}^*\|^2 &\leq \sum_{t=0}^f \left[ \frac{\|\mathbf{p}_0 - \mathbf{p}^*\|^2}{\nu\lambda\alpha t} + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} \alpha \right] \\ &\leq \frac{\|\mathbf{p}_0 - \mathbf{p}^*\|^2}{\nu\lambda\alpha} \log(f) + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f\alpha \\ &\leq \frac{\log(f)}{\nu\lambda f^{2/3}} \cdot \frac{1}{\alpha} + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f\alpha.\end{aligned}\quad (47)$$

Similar to (42) in Theorem B.9, we derive the positive projection for binding terms as

$$\begin{aligned}\mathbb{E}\left\|\left(\sum_{t=0}^f (\mathbf{d} - \mathbf{a}_t x_t)\right)^{B^+}\right\| &\leq \mathbb{E}\left\|\sum_{t=0}^f (\mathbf{d} - \mathbf{a}_t x_t)\right\| \\ &\leq \frac{1}{\alpha} \mathbb{E}\|\mathbf{p}_0 - \mathbf{p}_{f+1}\| \\ &\leq \frac{1}{\alpha} \mathbb{E}[\|\mathbf{p}_0 - \mathbf{p}^*\| + \|\mathbf{p}_{f+1} - \mathbf{p}^*\|] \\ &\leq \frac{1}{f^{1/3}} \cdot \frac{1}{\alpha} + \frac{C_w}{f^{1/3}\sqrt{f}} \cdot \frac{1}{\alpha\sqrt{\alpha}} + C_w \cdot \frac{1}{\sqrt{\alpha}} + \frac{C_w}{T^2} \cdot \frac{1}{\alpha}.\end{aligned}\quad (48)$$

Therefore, combining the results in (46), (47), and (48) together, we have

$$\begin{aligned}\Delta_{T_1 \text{ warm}} &\leq \frac{\log(f)}{\nu\lambda f^{2/3}} \cdot \frac{1}{\alpha} + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f\alpha \\ &\quad + \frac{2}{f^{1/3}} \cdot \frac{1}{\alpha} + \frac{2C_w}{f^{1/3}\sqrt{f}} \cdot \frac{1}{\alpha\sqrt{\alpha}} + 2C_w \cdot \frac{1}{\sqrt{\alpha}} + \frac{2C_w}{T^2} \cdot \frac{1}{\alpha} \\ &\leq \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f\alpha + \left( \frac{\log(f)}{\nu\lambda f^{2/3}} + \frac{2}{f^{1/3}} + \frac{2C_w}{T^2} \right) \cdot \frac{1}{\alpha} + \frac{2C_w}{f^{1/3}\sqrt{f}} \cdot \frac{1}{\alpha\sqrt{\alpha}} + 2C_w \cdot \frac{1}{\sqrt{\alpha}}.\end{aligned}\quad (49)$$

Thus, taking the optimal learning rate  $\alpha = f^{-2/3}$ , we obtain the regret for the first batch with a warm start as:

$$\Delta_{T_1 \text{ warm}} \leq \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f^{1/3} + \left( \frac{\log(f)}{\nu\lambda f^{2/3}} + \frac{2}{f^{1/3}} + \frac{2C_w}{T^2} \right) \cdot f^{2/3} + \frac{2C_w}{f^{1/3}\sqrt{f}} \cdot f + 2C_w \cdot f^{1/3}$$

$$\begin{aligned}
 &\leq \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} f^{1/3} + \frac{\log(f)}{\nu\lambda} + 2f^{1/3} + \frac{2C_w}{T} + 2C_w f^{1/6} + 2C_w \cdot f^{1/3} \\
 &\leq \left( \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} + 4C_w \right) \cdot f^{1/3} + \log(f).
 \end{aligned} \tag{50}$$

This completes the proof.

## C. Main Results

In this section, we demonstrate the proof for all theoretical results that we proposed in the main body of the paper. We instate Assumption 2.1 and 2.2 are satisfied.

### C.1. Proof of Theorem 3.1.

By Theorem B.3, we obtain

$$\begin{aligned}
 \Delta_T &= \mathbb{E}[r(\mathbf{x}) + v(\mathbf{x})] \\
 &\leq \|\mathbf{p}^*\| \cdot \mathbb{E} \left\| \left( \mathbf{b} - \sum_{t=1}^T \mathbf{a}_t x_t \right)^{B^+} \right\| + \frac{\mu \bar{a}^2}{2} \sum_{t=1}^T \mathbb{E} \|\mathbf{p}_t - \mathbf{p}^*\|^2 + \mathbb{E} \left\| \left( \sum_{t=1}^T \mathbf{a}_t x_t - \mathbf{b} \right)^+ \right\| \\
 &\leq \|\mathbf{p}^*\| \cdot \mathbb{E} [\|(\mathbf{b} - \mathbf{A}\mathbf{x})^{B^+}\|] + \mu \bar{a}^2 \sum_{t=1}^T \mathbb{E} [\|\mathbf{p}_t - \mathbf{p}^*\|^2] + \mathbb{E} [\|(\mathbf{A}\mathbf{x} - \mathbf{b})^+\|].
 \end{aligned}$$

This completes the proof.

### C.2. Proof of Theorem 3.2.

Combining the results of Theorem B.6, Theorem B.7, and Theorem B.9, for some constant  $C_{reg} > 0$ , we obtain:

$$\begin{aligned}
 \Delta_T &= \Delta_{T_1} + \Delta_{T_2} + \Delta_{T_3} \\
 &\leq \left( \frac{m(\bar{a} + \bar{d})^2}{2} + \frac{2\bar{r} + m(\bar{a} + \bar{d})^2}{\underline{d}} + m(\bar{a} + \bar{d}) \right) \sqrt{f} \\
 &\quad + \left( C_{lp} + \frac{mC}{\nu^2 \lambda^2} + (\bar{a} + \bar{d} + \delta) \frac{C}{\delta^2} \right) \log(k) \\
 &\quad + \left( \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} + \frac{2C_{lp}}{\sqrt{k}} + \frac{C_{lp} \log(f)}{\nu\lambda k} + \frac{4C_b}{\sqrt{k}} + 2C_b \right) \cdot f^{1/3}.
 \end{aligned} \tag{51}$$

Therefore, we achieve the worst-case regret of:

$$\Delta_T = \mathcal{O}(\log(k) + \sqrt{f} + f^{1/3}). \tag{52}$$

For special cases,

1. If we use LP-based method on the first batch, then we will have  $\Delta_{T_1} \leq \log(f)$ . Total performance follows

$$\Delta_T \leq \log(\max\{f, k\}) + f^{1/3} \leq \log(\sqrt{T}) + f^{1/3} \in \mathcal{O}(\log(T) + f^{1/3}). \tag{53}$$

2. If we have a warm start for the first batch, with the initialization  $\|\mathbf{p}_0 - \mathbf{p}^*\| \leq f^{-1/3}$ , then the first batch regret achieves  $\Delta_{T_1} \in \mathcal{O}(f^{1/3} + \log(f))$  by Lemma B.10. Thus, total performance follows

$$\Delta_T \leq \log(\max\{f, k\}) + 2f^{1/3} \in \mathcal{O}(\log(T) + f^{1/3}). \tag{54}$$

This completes our proof.

## D. Auxiliary Results

In this section, we provide auxiliary results to support the proof in the previous three sections. These lemmas focus on pure mathematical derivations.

### D.1. Technical Support for LP-based Analysis

**Lemma D.1.** Denote  $d_{i,(j+1)f}, d_{i,jf}$  as the average consumption of  $i$ -th type resource at time  $(j+1)f$  and  $jf$ . There exists a constant  $C_{rec}$  depending on  $\bar{d}, \bar{a}, m, \nu, \lambda, \mu$ , and  $C_{lp}$  such that:

$$\mathbb{E}(d_{i,(j+1)f} - d_i)^2 \leq \mathbb{E}(d_{i,jf} - d_i)^2 + \frac{C_{rec}}{(k-j-1)^2 f} + \frac{4\mu\bar{a}^2 \sqrt{C_{lp}}}{(k-j-1)\sqrt{(j+1)f}} \sqrt{\mathbb{E}[(d_{i,jf} - d_i)^2]}.$$

**Lemma D.2.** With the recursion relation in (26), there exists a constant  $C > 0$  depending on  $\bar{d}, \bar{a}, m, \nu, \lambda, \mu$ , and  $C_{lp}$  such that the summation of the total deviation of  $\mathbf{d}_t$  with the original  $\mathbf{d}$  satisfies:

$$\sum_{j=1}^k \mathbb{E}[(d_{i,jf} - d_i)^2] \leq \frac{C}{f} \log(k).$$

**Proof of Lemma D.1.** We analyze each term in (25). The key technique we use is to take conditional expectations and simplify the double summations.

(a) Term 1.

$$\begin{aligned} & \mathbb{E} \left[ \frac{(\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*))^2}{(T - (j+1)f)^2} \right] \\ &= \frac{1}{(T - (j+1)f)^2} \mathbb{E} \left[ \sum_{\ell=jf+1}^{(j+1)f} \sum_{s=jf+1}^{(j+1)f} \mathbb{E}[(d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*))(d_{i,jf} - a_{i,s} I(r_s > \mathbf{a}_s^\top \mathbf{p}_{jf}^*)) \mid \mathbf{d}_{jf}] \right] \\ &= \frac{1}{(T - (j+1)f)^2} \mathbb{E} \left[ \sum_{\ell=jf+1}^{(j+1)f} \mathbb{E}[(d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*))^2 \mid \mathbf{d}_{jf}] \right] \\ &\quad + \frac{1}{(T - (j+1)f)^2} \mathbb{E} \left[ \sum_{\ell \neq j, \ell=jf+1}^{(j+1)f} \mathbb{E}[(d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*))(d_{i,jf} - a_{i,j} I(r_j > \mathbf{a}_j^\top \mathbf{p}_{jf}^*)) \mid \mathbf{d}_{jf}] \right] \\ &\leq \frac{f(\bar{a} + \bar{d})^2}{(T - (j+1)f)^2} + \frac{\mathbb{E} \left[ \sum_{\ell \neq j, \ell=jf+1}^{(j+1)f} \mathbb{E}[d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) \mid \mathbf{d}_{jf}] \mathbb{E}[d_{i,jf} - a_{i,j} I(r_j > \mathbf{a}_j^\top \mathbf{p}_{jf}^*) \mid \mathbf{d}_{jf}] \right]}{(T - (j+1)f)^2} \\ &= \frac{f(\bar{a} + \bar{d})^2}{(T - (j+1)f)^2} \quad (\text{since } \mathbb{E}[d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) \mid \mathbf{d}_{jf}] = 0 \text{ for binding terms}) \\ &\leq \frac{(\bar{a} + \bar{d})^2}{(k-j-1)^2 f}. \quad (\text{since } \lfloor T \rfloor = k \cdot f) \end{aligned}$$

(b) Term 2.

$$\begin{aligned} & \mathbb{E} \left[ \frac{(\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}^*))^2}{(T - (j+1)f)^2} \right] \\ &= \frac{\mathbb{E} \left[ \sum_{\ell,s=jf+1}^{j(k+1)f} \mathbb{E}[(a_{i,\ell}(I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}^*))) \cdot (a_{i,s}(I(r_s > \mathbf{a}_s^\top \mathbf{p}_{jf}^*) - I(r_s > \mathbf{a}_s^\top \mathbf{p}_{(j+1)f}^*))) \mid \mathbf{d}_{jf}] \right]}{(T - (j+1)f)^2}. \end{aligned}$$

When  $\ell = s$ ,

$$\mathbb{E} \left[ (a_{i,\ell}(I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) - I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}))^2 \mid \mathbf{d}_{jf}) \right] \leq \bar{a}^2.$$

When  $\ell \neq s$ , by Assumption 2.2,

$$\begin{aligned} & \mathbb{E} \left[ (a_{i,\ell}(I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) - I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}))) (a_{i,s}(I(r_s > \mathbf{a}_s^\top \mathbf{p}_{jff}^*) - I(r_s > \mathbf{a}_s^\top \mathbf{p}_{(j+1)f}))) \mid \mathbf{d}_{jf} \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ (a_{i,\ell}(I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) - I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}))) (a_{i,s}(I(r_s > \mathbf{a}_s^\top \mathbf{p}_{jff}^*) - I(r_s > \mathbf{a}_s^\top \mathbf{p}_{(j+1)f}))) \mid \mathbf{a}_i, \mathbf{a}_s \right] \mid \mathbf{d}_{jf} \right] \\ &= \mathbb{E} \left[ a_{i,\ell} \mathbb{E} \left[ I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) - I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}) \mid \mathbf{a}_\ell \right] \cdot a_{i,s} \mathbb{E} \left[ I(r_s > \mathbf{a}_s^\top \mathbf{p}_{jff}^*) - I(r_s > \mathbf{a}_s^\top \mathbf{p}_{(j+1)f}) \mid \mathbf{a}_s \right] \mid \mathbf{d}_{jf} \right] \\ &= \mathbb{E} \left[ a_{i,\ell} (P(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^* \mid \mathbf{a}_\ell) - P(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f} \mid \mathbf{a}_\ell)) \cdot a_{i,s} (P(r_s > \mathbf{a}_s^\top \mathbf{p}_{jff}^* \mid \mathbf{a}_s) - P(r_s > \mathbf{a}_s^\top \mathbf{p}_{(j+1)f} \mid \mathbf{a}_s)) \mid \mathbf{d}_{jf} \right] \\ &\leq \mathbb{E} \left[ \mu a_{i,\ell} \mathbf{a}_\ell^\top (\mathbf{p}_{(j+1)f} - \mathbf{p}_{jff}^*) \cdot \mu a_{i,s} \mathbf{a}_s^\top (\mathbf{p}_{(j+1)f} - \mathbf{p}_{jff}^*) \mid \mathbf{d}_{jf} \right] \quad (\text{using Assumption here}) \\ &\leq \mu^2 \bar{a}^4 \mathbb{E} [(\mathbf{p}_{(j+1)f} - \mathbf{p}_{jff}^*)^2 \mid \mathbf{d}_{jf}]. \end{aligned}$$

Combining these two cases together and by the convergence of LP-based method, we obtain the bound for Term 2 as

$$\begin{aligned} & \mathbb{E} \left[ \frac{(\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}))^2}{(T - (j+1)f)^2} \right] \\ &\leq \frac{1}{(T - (j+1)f)^2} \mathbb{E} \left[ \sum_{\ell=jf+1}^{(j+1)f} \bar{a}^2 + \sum_{\ell \neq s, \ell=jf+1}^{(j+1)f} \mu^2 \bar{a}^4 \mathbb{E}[(\mathbf{p}_{(j+1)f} - \mathbf{p}_{jff}^*)^2 \mid \mathbf{d}_{jf}] \mid \mathbf{d}_{jf} \right] \\ &\leq \frac{1}{(T - (j+1)f)^2} \left( f \bar{a}^2 + f^2 \mu^2 \bar{a}^4 \mathbb{E}[(\mathbf{p}_{(j+1)f} - \mathbf{p}_{jff}^*)^2 \mid \mathbf{d}_{jf}] \right) \\ &\leq \frac{1}{(T - (j+1)f)^2} \left( f \bar{a}^2 + f^2 \mu^2 \bar{a}^4 \frac{C_{lp}}{jf} \right) \quad (\text{by Lemma A.4}) \\ &\leq \frac{\bar{a}^2 + \frac{1}{j} \mu^2 \bar{a}^4 C_{lp}}{(k - j - 1)^2 f}. \end{aligned}$$

(c) Term 3.

$$\begin{aligned} & 2\mathbb{E} \left[ (d_{i,jf} - d_i) \left( \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*)}{T - (j+1)f} \right) \right] \\ &= \frac{2}{T - (j+1)f} \mathbb{E} \left[ \sum_{\ell=jf+1}^{(j+1)f} \mathbb{E} [d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) \mid \mathbf{d}_{jf}] \cdot (d_{i,jf} - d_i) \right] \\ &= 0. \quad (\text{by the definition of binding terms}) \end{aligned}$$

(d) Term 4.

$$\begin{aligned} & 2\mathbb{E} \left[ (d_{i,jf} - d_i) \left( \frac{\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f} \right) \right] \\ &= \frac{2}{T - (j+1)f} \mathbb{E} \left[ \sum_{\ell=jf+1}^{(j+1)f} (d_{i,jf} - d_i) a_{i,\ell} \mathbb{E} [I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jff}^*) - I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f}) \mid \mathbf{a}_\ell] \right] \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{2}{T - (j+1)f} \sum_{\ell=jf+1}^{(j+1)f} \mathbb{E} \left[ (d_{i,jf} - d_i) a_{i,\ell} \left[ P(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^* \mid \mathbf{a}_\ell) - P(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f} \mid \mathbf{a}_\ell) \right] \right] \\
 &\leq \frac{2}{T - (j+1)f} \sum_{\ell=jf+1}^{(j+1)f} \mathbb{E} \left[ (d_{i,jf} - d_i) \mu a_{i,\ell} \mathbf{a}_\ell^\top (\mathbf{p}_{(j+1)f} - \mathbf{p}_{jf}^*) \right] \quad (\text{by Assumption 2.2}) \\
 &\leq \frac{2\mu\bar{a}^2}{T - (j+1)f} \sum_{\ell=jf+1}^{(j+1)f} \sqrt{\mathbb{E}[(d_{i,jf} - d_i)^2]} \cdot \sqrt{\mathbb{E}[(\mathbf{p}_{(j+1)f} - \mathbf{p}_{jf}^*)^2]} \quad (\text{by Cauchy's inequality}) \\
 &\leq \frac{2\mu\bar{a}^2}{T - (j+1)f} \cdot \frac{\sqrt{C_{lp}}f}{\sqrt{j}f} \sqrt{\mathbb{E}[(d_{i,jf} - d_i)^2]} \\
 &= \frac{4\mu\bar{a}^2 \sqrt{C_{lp}}}{(k-j-1)\sqrt{(j+1)f}} \sqrt{\mathbb{E}[(d_{i,jf} - d_i)^2]}.
 \end{aligned}$$

(e) Term 5.

$$\begin{aligned}
 &2\mathbb{E} \left( \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*)}{T - (j+1)f} \cdot \frac{\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f} I(jf < \tau) \right) \\
 &\leq 2\sqrt{\mathbb{E} \left( \frac{\sum_{\ell=jf+1}^{(j+1)f} d_{i,jf} - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*)}{T - (j+1)f} \right)^2} \sqrt{\mathbb{E} \left( \frac{\sum_{\ell=jf+1}^{(j+1)f} a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{jf}^*) - a_{i,\ell} I(r_\ell > \mathbf{a}_\ell^\top \mathbf{p}_{(j+1)f})}{T - (j+1)f} \right)^2} \\
 &\leq 2\sqrt{\frac{(\bar{a} + \bar{d})^2}{(k-j-1)^2 f}} \cdot \sqrt{\frac{\bar{a}^2 + \frac{1}{j}\mu^2\bar{a}^4 C_{lp}}{(k-j-1)^2 f}} \quad (\text{by results of Term 1 and Term 2}) \\
 &= \frac{2\bar{a}(\bar{a} + \bar{d})\sqrt{1 + \frac{1}{j}\mu^2\bar{a}^2 C_{lp}}}{(k-j-1)^2 f}.
 \end{aligned}$$

Combining all the terms, we obtain the upper bound as:

$$\begin{aligned}
 \mathbb{E}[(d_{i,(j+1)f} - d_i)^2] &\leq \mathbb{E}(d_{i,jf} - d_i)^2 + \frac{(\bar{a} + \bar{d})^2}{(k-j-1)^2 f} + \frac{\bar{a}^2 + \frac{1}{j}\mu^2\bar{a}^4 C_{lp}}{(k-j-1)^2 f} \\
 &\quad + \frac{2\bar{a}(\bar{a} + \bar{d})\sqrt{1 + \frac{1}{j}\mu^2\bar{a}^2 C_{lp}}}{(k-j-1)^2 f} + \frac{4\mu\bar{a}^2 \sqrt{C_{lp}}}{(k-j-1)\sqrt{(j+1)f}} \sqrt{\mathbb{E}[(d_{i,jf} - d_i)^2]}.
 \end{aligned}$$

Taking  $C_{rec} = (\bar{a} + \bar{d})^2 + \bar{a}^2 + \mu^2\bar{a}^4 C_{lp} + 2\bar{a}(\bar{a} + \bar{d})\sqrt{1 + \mu^2\bar{a}^2 C_{lp}}$  completes the proof.

**Proof of Lemma D.2.** We consider a general sequence  $\{z_j\}_{j=1}^k$  with

$$z_{j+1} \leq z_j + \frac{R}{(k-j-1)^2 f} + \frac{\sqrt{R}\sqrt{z_j}}{(k-j-1)\sqrt{(j+1)f}}$$

where  $R > 0$  is a constant.

Taking sum on both sides of the inequality and re-arranging, we have

$$\begin{aligned}
 \sum_{j=1}^k (k-j+1)(z_{j+1} - z_j) &\leq \sum_{j=1}^k \frac{16R}{(k-j+1)f} + \sqrt{\frac{16R}{f}} \sum_{j=1}^k \frac{\sqrt{z_j}}{\sqrt{j+1}} \\
 &\leq \frac{16R}{f} \log(k) + \sqrt{\frac{16R}{f}} \sum_{j=1}^k \frac{\sqrt{z_j}}{\sqrt{j+1}}.
 \end{aligned}$$



Noticing that  $\sum_{j=1}^k (k-j+1)(z_{j+1} - z_j) = \sum_{j=1}^k z_j$ , we have

$$\begin{aligned} \frac{16R}{f} \left( \sum_{j=1}^k \frac{1}{j+1} \right) \cdot \left( \sum_{j=1}^k z_j \right) &\geq \left( \sqrt{\frac{16R}{f}} \sum_{j=1}^k \frac{\sqrt{z_j}}{\sqrt{j+1}} \right)^2 \\ &\geq \left( \sum_{j=1}^k z_j - \frac{4R}{f} \log(k) \right)^2. \end{aligned}$$

We treat  $\sum_{j=1}^k z_j$  as the variable, then solve and get

$$\sum_{j=1}^k z_j \leq \frac{12R}{f} \log(k).$$

With this result, consider our recursion in (26), take the constant  $C = 12 \max\{C_{rec}, 16\mu^2\bar{a}^4 C_{lp}\} > 0$ , we obtain:

$$\sum_{j=1}^k \mathbb{E} [(d_{i,jf} - d_i)^2] \leq \frac{C}{f} \log(k).$$

## D.2. Technical Support for First-order Analysis

**Lemma D.3.** *Following the updated rule of first-order method, we derive the last dual price satisfies:*

$$\mathbb{E} \|\mathbf{p}_{T+1} - \mathbf{p}^*\| \leq \frac{C_{lp}}{\sqrt{\nu\lambda(T-kf)kf}} \cdot \frac{1}{\sqrt{\alpha}} + \frac{m(\bar{a} + \bar{d})}{\sqrt{\nu\lambda}} \cdot \sqrt{\alpha} + \frac{\bar{p}}{T^2}.$$

**Proof of Lemma D.3.** According to Lemma A.6, take  $\alpha_t = \alpha < \frac{1}{\nu\lambda}$ , we derive the conditional expectation as

$$\begin{aligned} \mathbb{E} [\|\mathbf{p}_{T+1} - \mathbf{p}^*\|^2 \mid \mathbf{p}_{kf}] &\leq (1 - \nu\lambda\alpha) \mathbb{E} [\|\mathbf{p}_T - \mathbf{p}^*\|^2 \mid \mathbf{p}_{kf}] + \alpha^2 m(\bar{a} + \bar{d})^2 \\ &\leq (1 - \nu\lambda\alpha)^{T-kf} \|\mathbf{p}_{kf} - \mathbf{p}^*\|^2 + \sum_{j=0}^{T-kf-1} \alpha^2 m(\bar{a} + \bar{d})^2 (1 - \nu\lambda\alpha)^j \\ &\leq (1 - \nu\lambda\alpha)^{T-kf} \|\mathbf{p}_{kf} - \mathbf{p}^*\|^2 + \frac{\alpha^2 m(\bar{a} + \bar{d})^2}{\nu\lambda\alpha} \\ &\leq \frac{1}{\nu\lambda\alpha(T-kf)} \|\mathbf{p}_{kf} - \mathbf{p}^*\|^2 + \frac{\alpha m(\bar{a} + \bar{d})^2}{\nu\lambda} \end{aligned} \quad (55)$$

where we use the technique of  $\sum_{j=0}^{T-kf-1} \alpha^2 m(\bar{a} + \bar{d})^2 (1 - \nu\lambda\alpha)^j \leq \frac{1 - (1 - \nu\lambda\alpha)^{T-kf}}{\nu\lambda\alpha} \leq \frac{1}{\nu\lambda\alpha}$  and  $(1 - \nu\lambda\alpha)^{T-kf} \leq \frac{1}{1 + \nu\lambda\alpha(T-kf)} \leq \frac{1}{\nu\lambda\alpha(T-kf)}$ .

By LP-convergence result in Lemma A.4, we know  $\mathbb{E} \|\mathbf{p}_{kf} - \mathbf{p}^*\| \leq \frac{C_{lp}}{\sqrt{kf}}$ . By Proposition 3.3 in (Gao et al., 2024), we know the event  $E := \|\mathbf{p}_{kf} - \mathbf{p}^*\| \leq \frac{C_{lp}}{\sqrt{kf}}$  with probability  $\mathbb{P} \geq 1 - \frac{1}{T^4}$ . By Lemma B.2, we know  $\|\mathbf{p}_t\| \leq \bar{p}$ .

Thus, we have

$$\begin{aligned} \mathbb{E} \|\mathbf{p}_{T+1} - \mathbf{p}^*\|^2 &\leq \mathbb{E} [\|\mathbf{p}_{T+1} - \mathbf{p}^*\|^2 \mid E] \cdot \mathbb{P}(E) + \mathbb{E} [\|\mathbf{p}_{T+1} - \mathbf{p}^*\|^2 \mid \bar{E}] \cdot \mathbb{P}(\bar{E}) \\ &\leq \frac{C_{lp}}{\nu\lambda(T-kf)kf} \cdot \frac{1}{\alpha} + \frac{m(\bar{a} + \bar{d})^2}{\nu\lambda} \cdot \alpha + \frac{\bar{p}}{T^4}. \end{aligned} \quad (56)$$

Thus, as  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for any  $a, b > 0$ , by (56), we have

$$\mathbb{E}\|\mathbf{p}_{T+1} - \mathbf{p}^*\| \leq \frac{C_{lp}}{\sqrt{\nu\lambda}(T-kf)kf} \cdot \frac{1}{\sqrt{\alpha}} + \frac{m(\bar{a} + \bar{d})}{\sqrt{\nu\lambda}} \cdot \sqrt{\alpha} + \frac{\bar{p}}{T^2}.$$

This completes the proof.

## E. Supplementary Experiments

In this section, we provide more experiments to further evaluate the performance of our algorithms. We consider a more general and complex input distribution and include additional comparisons with recent methods.

### E.1. New Distribution

As an extension of Section 4.1, our goal is to evaluate the main algorithms with different choices of re-solving frequency  $f \in \{T^{1/3}, T^{1/2}, T^{2/3}\}$ . We consider a more complex distribution for reward and resource consumption requests and guarantee that Assumptions 2.1 and 2.2 are still satisfied.

Consider

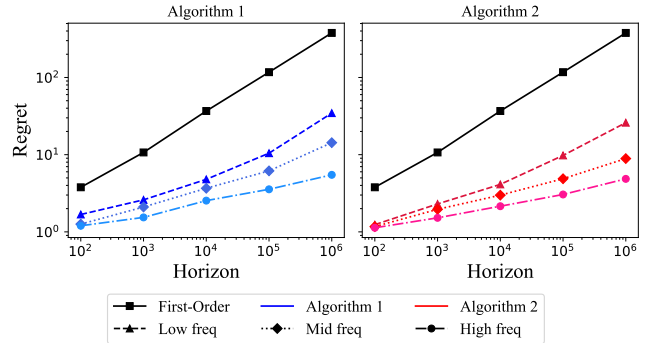
Input III:  $a_{it} \sim \min(1, \max(0, 1 + z))$ ,  $r_t \sim \text{Unif}[0, 1]$   
 where  $z \sim t(1)$  : Student's t-distribution with 1 degree of freedom.

We generate  $\{r_t, \mathbf{a}_t\}_{t=1}^T$  from Input III and keep other parameters  $T \in [10^2, 10^6]$  and  $d_i \sim \text{Uniform}[1/3, 2/3]$  the same in Section 4.1. We report the average result over 100 trials for each experiment and use the classic first-order method with  $\mathcal{O}(T^{1/2})$  regret (Algorithm 4) as a baseline.

Table 4: Algorithms under New Distribution.

	$T$	First-Order	Low freq	Mid freq	High freq
Algorithm 1	$10^2$	3.78	1.68	1.26	<b>1.20</b>
	$10^3$	10.69	2.60	2.10	<b>1.54</b>
	$10^4$	36.79	4.79	3.67	<b>2.54</b>
	$10^5$	117.13	10.45	6.17	<b>3.56</b>
	$10^6$	377.84	34.54	14.30	<b>5.48</b>
Algorithm 2	$10^2$	3.78	1.24	1.17	<b>1.13</b>
	$10^3$	10.69	2.30	1.96	<b>1.52</b>
	$10^4$	36.79	4.12	2.99	<b>2.15</b>
	$10^5$	117.13	9.78	4.88	<b>3.05</b>
	$10^6$	377.84	25.96	8.90	<b>4.86</b>

Figure 5: Regret for various re-solving frequencies.



Under the new distribution, our algorithms still exhibit a strong regret performance. As shown in Table 4 and Figure 5, we observe that regret decreases as the re-solving frequency increases. This trend holds consistently across both algorithms and is consistent with the guarantees of Theorem 3.2.

### E.2. More Comparison

Building on the comparison of baseline methods in Section 4.2, we evaluate our algorithms against recent works. As a reminder, our algorithms employ frequent LP-solving to learn online dual price under continuous support. Li et al. (2024) considers a similar problem using infrequent LP-solving to update the dual variable but under finite support.

To compare them, we adapt our algorithms to finite support and take the re-solving frequency  $f = T^{1/3}$  from Section 4.1; We take the best-performed parameters for the infrequent method where the number of customer types  $n = 50$  and  $\alpha = \beta = 0.95$  which control the solving times near start and end. We take the horizon over  $T \in [10^2, 10^6]$  and report the average result over 100 trials for each experiment. We still use the first-order method in Algorithm 4 as a baseline.

Table 5: Algorithms comparison.

$T$	Regret	Algorithm	Compute Time (s)
$10^3$	11.92	First-Order	0.002
	11.12	Infrequent LP-based	0.724
	<b>3.77</b>	Algorithm 1	<b>0.790</b>
	<b>3.69</b>	Algorithm 2	<b>0.782</b>
$10^4$	36.37	First-Order	0.01
	14.20	Infrequent LP-based	0.827
	<b>6.69</b>	Algorithm 1	<b>3.015</b>
	<b>6.35</b>	Algorithm 2	<b>3.112</b>
$10^5$	110.22	First-Order	0.109
	20.90	Infrequent LP-based	1.028
	<b>16.41</b>	Algorithm 1	<b>52.95</b>
	<b>10.84</b>	Algorithm 2	<b>52.39</b>
$10^6$	312.26	First-Order	1.106
	28.80	Infrequent LP-based	4.929
	<b>20.58</b>	Algorithm 1	<b>1261.0</b>
	<b>14.30</b>	Algorithm 2	<b>1305.4</b>

Figure 6: Regret for various algorithms.

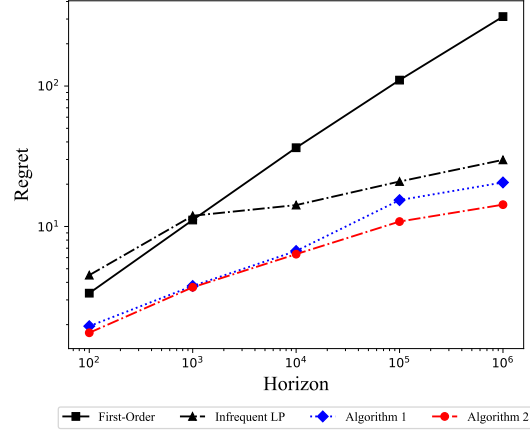


Table 5 and Figure 6 demonstrate the algorithm regret and computation time across different horizons. While the infrequent LP-based method has faster computation, our algorithms show a competitive performance in decision optimality and achieve lower regret under finite support.