# Efficient Online Learning with Memory via Frank-Wolfe Optimization: Algorithms with Bounded Dynamic Regret and Applications to Control

Hongyu Zhou [1]    Zirui Xu [1]    Vasileios Tzoumas [1]

## Abstract

Projection operations are a typical computation bottleneck in online learning. In this paper, we enable projection-free online learning within the framework of *Online Convex Optimization with Memory* (OCO-M) —OCO-M captures how the history of decisions affects the current outcome by allowing the online learning loss functions to depend on both current and past decisions. Particularly, we introduce the first projection-free meta-base learning algorithm with memory that minimizes dynamic regret, *i.e.*, that minimizes the suboptimality against *any* sequence of time-varying decisions. We are motivated by artificial intelligence applications where autonomous agents need to adapt to time-varying environments in real-time, accounting for how past decisions affect the present. Examples of such applications are: online control of dynamical systems; statistical arbitrage; and time series prediction. The algorithm builds on the Online Frank-Wolfe (**OFW**) and **Hedge** algorithms. We demonstrate how our algorithm can be applied to the online control of linear time-varying systems in the presence of unpredictable process noise. To this end, we develop a controller with memory and bounded dynamic regret against any optimal time-varying linear feedback control policy. We validate our algorithm in simulated scenarios of online control of linear time-invariant systems.

## 1. Introduction

Online Convex Optimization (OCO) (Shalev-Shwartz et al., 2012; Hazan et al., 2016) has found widespread application in statistics, information theory, and operation research (Cesa-Bianchi & Lugosi, 2006). OCO can be interpreted as a sequential game

between an optimizer and an adversary over $T$ time steps: at each time step $t = 1, \ldots, T$, first the optimizer chooses a decision $\mathbf{x}_t$ from a convex set $\mathcal{X}$; then, the adversary reveals a convex loss function $f_t$ and the optimizer suffers the loss $f_t(\mathbf{x}_t)$. The optimizer aims to minimize its cumulative loss, despite knowing each $f_t$ only after $\mathbf{x}_t$ has been already decided.

*Static regret* is the standard approach to measure the suboptimality of the optimizer's decisions $\mathbf{x}_1, \ldots, \mathbf{x}_T$. Particularly, given a decision $\mathbf{x} \in \mathcal{X}$ to compare $\mathbf{x}_1, \ldots, \mathbf{x}_T$ with, the static regret of $\mathbf{x}_1, \ldots, \mathbf{x}_T$ with respect to $\mathbf{v}$ is defined as follows (Hazan et al., 2016):

$$\text{S-Reg}_T = \sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{v}). \tag{1}$$

That is, when $\mathbf{v}$ minimizes $\sum_{t=1}^{T} f_t(\mathbf{v})$, then $\text{S-Reg}_T$ captures the suboptimality of $\mathbf{x}_1, \ldots, \mathbf{x}_T$ against the optimal *static* decision that would have been made in hindsight.

Algorithms that guarantee *static no-regret* have been widely adopted in applications pertained to recommendation systems, communication-channel allocation, and action prediction (Cesa-Bianchi & Lugosi, 2006).[1]

But the application of such algorithms to complex artificial intelligence tasks such as *online control under unpredictable disturbances* (Shi et al., 2019) and *collaborative multi-robot motion planning* (Xu et al., 2023) is hindered by three main technological challenges:

- **Challenge I: Dynamic Environments.** Complex tasks such as the above require decisions that adapt to changing environments. For example, *target tracking with multiple robots* requires the robots to continuously change their position to track moving targets (Xu et al., 2023). Therefore, measuring performance against a static (optimal) decision per the static regret in eq. (1) is insufficient. Instead, we need to measure performance against *time-varying* (optimal) decisions.

- **Challenge II: Past Decisions Affect the Present.** In complex tasks such as the aforementioned, past decisions

[1]Department of Aerospace Engineering, University of Michigan, Ann Arbor. Correspondence to: Hongyu Zhou <zhouhy@umich.edu>.

*Preliminary work.*

---

[1]An algorithm has *static no-regret* when $\text{S-Reg}_T/T$ tends to 0 when $T$ tends to $+\infty$, implying $f_t(\mathbf{x}_t)$ tends to $f_t(\mathbf{v})$ for $t$ large.

often affect the present outcome. Therefore, the OCO framework we discussed above, where each loss function $f_t$ depends on the most recent decision $\mathbf{x}_t$ only, fails to capture the effect of earlier decisions to the present. Instead, we need an OCO framework with memory, where each loss function $f_t$ depends on $\mathbf{x}_t$ as well as on the past $\mathbf{x}_{t-m}, \ldots, \mathbf{x}_{t-1}$, for some $m \geq 0$.

- **Challenge III: Fast Decision-Making.** Complex control tasks often require decisions to be made fast. For example, such is the case for the effective *online control of quadrotors against wind disturbances* (Romero et al., 2022). But the current OCO algorithms typically rely on projection operations which can be computationally expensive since they require solving quadratic programs (Kalhan et al., 2021). Instead, we need fast OCO algorithms that are inevitably projection-free.

All in all, the above challenges give rise to the need below:

**Need.** *We need online learning algorithms for OCO with Memory (OCO-M) that are projection-free and guarantee near-optimal decisions in dynamic environments. The decisions' near-optimality may be captured by bounding their suboptimality with respect to optimal decisions that adapt to the changing environment knowing its future evolution, i.e., by bounding* dynamic regret.

*Dynamic regret for the classical OCO without memory* is defined as follows (Zinkevich, 2003): given a time-varying comparator sequence $\mathbf{v}_1, \ldots, \mathbf{v}_T$, then[2]

$$\text{D-Reg}_T = \sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{v}_t). \qquad (2)$$

Dynamic regret contrasts static regret: static regret compares $(\mathbf{x}_1, \ldots, \mathbf{x}_T)$ against a merely static $\mathbf{v}$. Thus, when $\mathbf{v}_1, \ldots, \mathbf{v}_T$ minimize $\sum_{t=1}^{T} f_t(\mathbf{v}_t)$, then $\text{D-Reg}_T$ captures the suboptimality of $\mathbf{x}_1, \ldots, \mathbf{x}_T$ against the optimal *time-varying* decisions that would have been made in hindsight. Hence, dynamic regret bounds are typically larger than static regret bounds, depending on terms that capture the change of the environment. Such terms are *loss variation* $V_T$, *gradient variation* $D_T$, and *path length* $C_T$:[3]

$$V_T \triangleq \sum_{t=1}^{T} \sup_{\mathbf{x} \in \mathcal{X}} |f_t(\mathbf{x}) - f_{t-1}(\mathbf{x})|, \qquad (3)$$

$$D_T \triangleq \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2, \qquad (4)$$

$$C_T \triangleq \sum_{t=1}^{T} \|\mathbf{v}_t - \mathbf{v}_{t-1}\|_2. \qquad (5)$$

*Dynamic regret for OCO-M with memory* $m$, where the loss function at each time step $t$ takes the form $f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t) : \mathcal{X}^{m+1} \mapsto \mathbb{R}$, is defined as follows:

$$\text{Regret}_T^D = \sum_{t=1}^{T} f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{v}_{t-m}, \ldots, \mathbf{v}_t), \qquad (6)$$

where it is assumed that $\mathbf{x}_{t-m} = \mathbf{0}$ for $t \leq m$.

**Contributions.** We aim to address the Need by means of the following contributions:

- *Algorithmic Contributions*: We introduce the first projection-free algorithm for OCO-M with bounded dynamic regret (Sections 4 and 5) —the regret bound is presented in Table 1. The algorithm builds on the projection-free algorithms **Hedge** (Freund & Schapire, 1997) and **OFW** (Hazan & Kale, 2012; Kalhan et al., 2021).

  We apply our algorithm to the online control of linear time-varying systems in the presence of unpredictable noise (Section 6). We thus introduce a projection-free controller with memory and bounded dynamic regret against any optimal time-varying linear feedback control gains. Particularly, our comparator class of optimal time-varying linear feedback control gains does <u>not</u> require the a priori knowledge of stabilizing control gains. Instead, the state-of-the-art OCO-M controller by (Zhao et al., 2022) requires a comparator class of optimal time-varying policies where an a priori knowledge of stabilizing control gains is necessary.

- *Technical Contributions*: To enable aforementioned algorithmic and regret bound contributions, we make the following technical innovations:

  - We analyze dynamic regret of the **OFW** algorithm (Section 4). The analysis enables the state-of-the-art bound in (Kalhan et al., 2021, Theorem 1) to hold true for any convex loss functions in the evaluation of **OFW**'s regret (see Table 1). Instead, (Kalhan et al., 2021, Theorem 1) holds true for smooth convex functions only.

  - We prove that the *Disturbance-Action Control* (DAC) policy (Agarwal et al., 2019) —widely used in online non-stochastic control to reduce the online control problem to OCO-M (Agarwal et al., 2019; Gradu et al., 2020; Hazan et al., 2020)— is able to approximate time-varying linear feedback controllers (Proposition 2 in Appendix D.5). Previous results have established that a DAC policy can approximate time-<u>in</u>variant linear feedback controllers only (Agarwal et al., 2019), instead of a time-varying controllers.

---

[2]A related measure to dynamic regret is *adaptive regret* (Hazan & Seshadhri, 2007). Adaptive regret captures the worst-case static regret on any contiguous time interval. (Zhang, 2020) studies the relation of dynamic regret to adaptive regret.

[3]Obtaining a no-regret algorithm hence requires the growth of the metrics in eqs. (3) to (5) to be sublinear (Besbes et al., 2015; Mokhtari et al., 2016; Kalhan et al., 2021). $V_T$ and $D_T$ are small when the loss function and decisions change slowly.

*Table 1.* **Comparison of related work and our work on contributed algorithms with bounded dynamic regret bounds for Online Convex Optimization**. GO denotes the number of gradient oracle calls per iteration of the respective algorithm.

| Reference | Loss function | Projection-free | Memory | GO | Regret Rate |
|---|---|---|---|---|---|
| (Zinkevich, 2003) | Convex | No | No | $\mathcal{O}(1)$ | $\mathcal{O}(\sqrt{T}(1+C_T))$ |
| (Jadbabaie et al., 2015) | Convex smooth | No | No | $\mathcal{O}(1)$ | $\mathcal{O}\left(\sqrt{(1+D_T)}+\min\left\{\sqrt{(1+D_T)C_T},(1+D_T)^{\frac{1}{3}}T^{\frac{1}{3}}V_T^{\frac{1}{3}}\right\}\right)$ |
| (Mokhtari et al., 2016) | Strongly convex | No | No | $\mathcal{O}(1)$ | $\mathcal{O}(1+C_T)$ |
| (Yang et al., 2016) | Convex smooth | No | No | $\mathcal{O}(1)$ | $\mathcal{O}(C_T)$ |
| (Zhang et al., 2018) | Convex | No | No | $\mathcal{O}(1)$ | $\mathcal{O}(\sqrt{T(1+C_T)})$ |
| (Kalhan et al., 2021) | Convex smooth | **Yes** | No | $\mathcal{O}(1)$ | $\mathcal{O}\left(\sqrt{T}(1+V_T+\sqrt{D_T})\right)$ |
| Ours (Theorem 1 and Theorem 4) | Convex | **Yes** | No | $\mathcal{O}(1)$ | $\mathcal{O}\left(\sqrt{T}(1+V_T+\sqrt{D_T}),\mathcal{O}\left(\sqrt{T(V_T+D_T)}\right)\right)$ |
| (Zhao et al., 2022) | Convex | No | **Yes** | $\mathcal{O}(1)$ | $\mathcal{O}(\sqrt{T(1+C_T)})$ |
| Ours (Theorem 2) | Convex | **Yes** | **Yes** | $\mathcal{O}(1)$ | $\mathcal{O}(\sqrt{T(1+V_T+D_T+C_T)})$ |

**Numerical Evaluations.** We validate our algorithm in simulated scenarios of online control of linear time-invariant systems (Appendix E). We compare our algorithm with **OGD** (Zinkevich, 2003), **Ader** (Zhang et al., 2018), and **Scream** (Zhao et al., 2022) algorithms. Our algorithm is observed 3 times faster than the state-of-the-art OCO-M algorithm **Scream** (Zhao et al., 2022) as system dimension increases, and achieves comparable or superior loss performance over all compared algorithms.

## 2. Related Work

We review the literature by first reviewing *OCO without Memory* and *OCO with Memory*; then, we review *Online Learning for Control via OCO with Memory*.

**OCO without Memory.** The *OCO without Memory* literature is vast (Hazan et al., 2016). We here focus on algorithms that guarantee bounded dynamic regret; a representative subset is presented in Table 1.

(Zhang et al., 2018) prove that the optimal dynamic regret for OCO without Memory is $\Omega\left(\sqrt{T(1+C_T)}\right)$, and provide an algorithm matching this bound. The algorithm is based on *Online Gradient Descent* (**OGD**), which is a projection-based algorithm: at each time step $t$, **OGD** chooses a decision $\mathbf{x}_t$ by first computing an intermediate decision $\mathbf{x}'_t = \mathbf{x}_{t-1} - \eta\nabla f_{t-1}(\mathbf{x}_{t-1})$ —given the previous decision $\mathbf{x}_{t-1}$, the gradient of the previously revealed loss $f_{t-1}(\mathbf{x}_{t-1})$, and a step size $\eta > 0$— and then projects $\mathbf{x}'_t$ back to the feasible convex set $\mathcal{X}$ to output the final decision $\mathbf{x}_t$. This projection operation is often computationally expensive since it requires solving a quadratic program (Rockafellar, 1976). When the projection operation is indeed computationally expensive, the *Online Frank-Wolfe* (**OFW**) algorithm is employed as a projection-free alternative (Frank & Wolfe, 1956; Hazan & Kale, 2012): **OFW** seeks a feasible descent direction by solving the linear program $\mathbf{x}'_{t-1} = \arg\min_{\mathbf{x}\in\mathcal{X}}\langle\nabla f_{t-1}(\mathbf{x}_{t-1}),\mathbf{x}\rangle$ and then updating $\mathbf{x}_t = (1-\eta)\mathbf{x}_{t-1} + \eta\mathbf{x}'_{t-1}$. (Kalhan et al., 2021) generalize the **OFW** method to OCO without Memory to achieve a bounded dynamic regret and **OFW** has been ob-

served 20 times faster than **OGD** (Kalhan et al., 2021).[4]

**OCO with Memory.** (Zhao et al., 2022) prove that the optimal dynamic regret for OCO-M is $\Omega(\sqrt{T(1+C_T)})$, and provide an algorithm matches thing bound based on **OGD**. Earlier works have provided static regret bounds for OCO-M, such as the bound $\mathcal{O}(T^{2/3})$ by (Weinberger & Ordentlich, 2002), and the bound $\mathcal{O}(\sqrt{T})$ by (Anava et al., 2015). We provide the first projection-free algorithm for OCO-M that also guarantees bounded dynamic regret.

**Online Learning for Control via OCO-M.** OCO-M has been recently applied to the control of linear dynamical systems in the presence of adversarial (non-stochastic) noise (Agarwal et al., 2019; Simchowitz et al., 2020; Shalev-Shwartz et al., 2012). The noise is adversarial in the sense that it may adapt to the system's evolution. Generally, the noise can evolve arbitrarily, subject to a given upper bound on its magnitude —the upper bound ensures problem feasibility. Thus, no stochastic model is assumed regarding the noise's evolution, in contrast to classical control that typically assumes Gaussian noise (Åström, 2012).

The current OCO-M algorithms for control prescribe control policies by optimizing linear feedback control gains. The algorithms rely on projection-based methods such as **OGD**, and guarantee bounded static regret (Agarwal et al., 2019; Hazan et al., 2020; Simchowitz et al., 2020; Li et al., 2021), adaptive regret (Gradu et al., 2020; Zhang et al., 2022), or dynamic regret (Zhao et al., 2022). Specifically, the said OCO-M regret bounds are against optimal static feedback control gains with the exception of the bound by (Zhao et al., 2022) which is against a class of optimal time-varying policies; however, the definition of this

---

[4]Additional examples of works utilizing **OFW** for OCO without Memory are: (Hazan & Kale, 2012; Jaggi, 2013; Garber & Hazan, 2015; Wan & Zhang, 2021; Kalhan et al., 2021; Kretzu & Garber, 2021; Garber & Kretzu, 2022; Wan et al., 2023). Examples of works utilizing **OGD** for OCO without Memory are: (Zinkevich, 2003; Jadbabaie et al., 2015; Mokhtari et al., 2016; Yang et al., 2016; Zhang et al., 2018; Chang & Shahrampour, 2020).

class requires an a priori knowledge of linear feedback control gains that ensure stability. We provide a projection-free controller with memory and bounded dynamic regret against any optimal time-varying linear feedback control policy without the need to specify to the optimal policy any stabilizing feedback control gains.

## 3. Problem Formulation

We formally define the problem of *Online Convex Optimization with Memory* (OCO-M) (Problem 1), along with standard convexity (but non-smoothness) assumptions.

**Problem 1** (Online Convex Optimization with Memory (OCO-M) (Weinberger & Ordentlich, 2002)). *There exist 2 players, an online optimizer and an adversary, who choose decisions sequentially over a time horizon $T$. At each time step $t = 1, \ldots, T$, the online optimizer chooses a decision $\mathbf{x}_t$ from a convex set $\mathcal{X}$; then, the adversary chooses a loss $f_t : \mathcal{X}^{m+1} \mapsto \mathbb{R}$ to penalize the optimizer's most recent $m + 1$ decisions. Particularly, the adversary reveals $f_t$ to the optimizer and the optimizer computes its loss $f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t)$, where $\mathbf{x}_{t-m}$ is $\mathbf{0}$ for $t \leq m$. The optimizer aims to minimize $\sum_{t=1}^{T} f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t)$.*

The challenge in solving OCO-M optimally, *i.e.*, in minimizing $\sum_{t=1}^{T} f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t)$, is that the optimizer gets to know $f_t$ only after $\mathbf{x}_t$ has been chosen, instead of before.

Despite the above challenge, our objective is to develop an efficient (projection-free) online algorithm for OCO-M that despite its efficiency still enjoys sublinear dynamic regret.

To achieve our objective, we adopt standard assumptions in online convex optimization (Hazan et al., 2016; Anava et al., 2015; Agarwal et al., 2019; Simchowitz et al., 2020; Zhang et al., 2018; Gradu et al., 2020; Zhao et al., 2022):

**Assumption 1** (Convex and Compact Bounded Domain, Containing the Origin). *The domain set $\mathcal{X}$ is convex and compact, contains the zero point, and has diameter $D$, where $D$ is a given non-negative number; i.e., $\mathbf{0} \in \mathcal{X}$, and $\|\mathbf{x} - \mathbf{y}\|_2 \leq D$ for all $\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{X}$.*

**Definition 1** (Unary Loss Function). *Given $f_t : \mathcal{X}^{m+1} \mapsto \mathbb{R}$, the unary loss function is the $\widetilde{f}_t(\mathbf{x}) \triangleq f_t(\mathbf{x}, \ldots, \mathbf{x})$.*

**Assumption 2** (Convex Loss). *The loss function $f_t : \mathcal{X}^{m+1} \mapsto \mathbb{R}$ is convex, i.e., the unary loss function $\widetilde{f}_t(\mathbf{x})$ is convex in $\mathbf{x}$, where $m$ is the memory length, and $\mathbf{x} \in \mathcal{X}$.*

**Assumption 3** (Bounded Loss). *The loss function $f_t$ takes values in $[a, a + c]$, where $a$ and $c$ are non-negative; i.e.,*

$$0 \leq a \leq f_t(\mathbf{x}_0, \ldots, \mathbf{x}_m) \leq a + c,$$

*for all $(\mathbf{x}_0, \ldots, \mathbf{x}_m) \in \mathcal{X}^{m+1}$ and $t \in \{1, \ldots, T\}$.*

**Assumption 4** (Coordinate-Wise Lipschitz). *The loss function $f_t$ is coordinate-wise $L$-Lipschitz, where $L$ is a given*

non-negative number; i.e.,

$$|f_t(\mathbf{x}_0, \ldots, \mathbf{x}_m) - f_t(\mathbf{y}_0, \ldots, \mathbf{y}_m)| \leq L \sum_{i=0}^{m} \|\mathbf{x}_i - \mathbf{y}_i\|_2,$$

*for all $(\mathbf{x}_0, \ldots, \mathbf{x}_m) \in \mathcal{X}^{m+1}$, and $(\mathbf{y}_0, \ldots, \mathbf{y}_m) \in \mathcal{X}^{m+1}$, and for all $t \in \{1, \ldots, T\}$.*

**Assumption 5** (Bounded Gradient). *The gradient norm of $\widetilde{f}_t$ is at most $G$, where $G$ is a given non-negative number; i.e., $\left\|\nabla \widetilde{f}_t(\mathbf{x})\right\|_2 \leq G$ for all $\mathbf{x} \in \mathcal{X}$ and $t \in \{1, \ldots, T\}$.*

## 4. Meta-OFW Algorithm for OCO-M

We present **Meta-OFW**, the first projection-free algorithm with bounded dynamic regret for OCO-M. **Meta-OFW** leverages as subroutine the Online Frank-Wolfe (**OFW**) algorithm. **OFW** is introduced by (Kalhan et al., 2021) for the OCO problem without memory.

We next first present the **OFW** algorithm (Section 4.1), and then present the **Meta-OFW** algorithm (Section 4.2).

### 4.1. The Online Frank-Wolfe (**OFW**) Algorithm for OCO without Memory

We present the **OFW** algorithm (Algorithm 1) along with its dynamic regret analysis (Theorem 1). Particularly, our analysis results in the same regret bound as **OFW**'s state of the art bound in (Kalhan et al., 2021, Theorem 1) but under Assumption 1 and Assumption 2 only. Instead, **OFW**'s bound in (Kalhan et al., 2021) holds true under the additional assumption of smooth loss functions.

**Theorem 1** (Dynamic Regret Bound of **OFW** for OCO with no memory). *Consider the OCO problem with no memory, i.e., Problem 1 with $m = 0$. Under Assumption 1 and Assumption 2, **OFW** achieves against any sequence of comparators $(\mathbf{v}_1, \ldots, \mathbf{v}_T) \in \mathcal{X}^T$ the dynamic regret*

$$\text{Regret}_T^D \leq \mathcal{O}\left(\frac{1 + V_T}{\eta} + \sqrt{T D_T}\right). \quad (7)$$

*Particularly, when $\eta$ is chosen such that $\eta = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$, then*

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T}\left(1 + V_T + \sqrt{D_T}\right)\right). \quad (8)$$

The **OFW** algorithm achieves Theorem 1 by executing the following projection-free steps (Algorithm 1): **OFW** first takes as input the time horizon $T$ and a constant step size $\eta$. Then, at each iteration $t = 1, \ldots, T$, **OFW** chooses an $\mathbf{x}_t$, after which the learner suffers a loss $f_t(\mathbf{x}_t)$ and evaluates the gradient $\nabla f_t(\mathbf{x}_t)$ (lines 3-4). Afterwards, **OFW** seeks a direction $\mathbf{x}_t'$ that is parallel to the gradient within the feasible set $\mathcal{X}$ by solving a linear program only once per iteration (line 5). Finally, the decision for next iteration is then updated by $\mathbf{x}_{t+1} = (1 - \eta)\mathbf{x}_t + \eta \mathbf{x}_t'$ (line 6).

**Algorithm 1** Online Frank-Wolfe Algorithm (**OFW**) (Kalhan et al., 2021).

**Input:** Time horizon $T$; step size $\eta$.
**Output:** Decision $\mathbf{x}_t$ at each time step $t = 1, \ldots, T$.

1: Initialize $\mathbf{x}_1 \in \mathcal{X}$;
2: **for** each time step $t = 1, \ldots, T$ **do**
3:     Suffer a loss $f_t(\mathbf{x}_t)$;
4:     Obtain gradient $\nabla f_t(\mathbf{x}_t)$;
5:     Compute $\mathbf{x}'_t = \arg\min_{\mathbf{x} \in \mathcal{X}} \langle \nabla f_t(\mathbf{x}_t), \mathbf{x} \rangle$;
6:     Update $\mathbf{x}_{t+1} = (1-\eta)\mathbf{x}_t + \eta\mathbf{x}'_t$;
7: **end for**

**Remark 1** (Efficiency due to only Projection-Free Operations). **OFW** *in Algorithm 1 is projection-free: it finds a descent direction within the feasible set via solving a linear program once per iteration (line 5). Instead, e.g.,* **OGD** *requires solving a quadratic program for projections (Zinkevich, 2003). Thus,* **OFW** *is more efficient when projections are costly. For example, (Kalhan et al., 2021) demonstrates that* **OFW** *is 20 times faster than* **OGD** *in matrix completion scenarios. In the numerical evaluations in this paper (Appendix E), over online non-stochastic control scenarios, we observe that the proposed* **OFW***-based algorithm is about 3 times faster than the* **OGD***-based algorithm (achieving comparable or superior loss performance).*

### 4.2. Meta-OFW Algorithm for OCO-M

We present **Meta-OFW** (Algorithm 2). To this end, we start with the intuition on how Algorithm 2's steps achieve a bounded dynamic regret (the rigorous dynamic regret analysis of **Meta-OFW** is given in Section 5).

Algorithm 2 utilizes multiple copies of the **OFW** algorithm as base-learners —each one with a different step size $\eta$— and the **Hedge** algorithm (Freund & Schapire, 1997) as a meta-learner. The multiple copies of **OFW** aim to cope with the a priori unknown loss variation $V_T$ via a trick reminiscent of the "doubling trick" (Shalev-Shwartz et al., 2012), *i.e.*, via covering the spectrum of step sizes such that there exist a step size that approximately minimizes eq. (7) as if $V_T$ was known; and **Hedge** fuses the decisions provided by the base-learners to output a final decision $\mathbf{x}_t$ at each step $t$.

We discuss in more detail the role of the base- and meta-learners in Remark 2 and Remark 3 below, respectively. To this end, we use the following notation and definitions:

- $\lambda \triangleq m^2 L$ is a regularizing constant;
- $N$ is the total number of the base-learners;
- $\mathcal{B}_i$ is the $i$-th base-learner running **OFW** with step size $\eta_i$ and output $\mathbf{x}_{t,i}$ at each iteration $t$, where $i \in \{1, \ldots, N\}$;
- $g_t(\mathbf{x}) \triangleq \langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{x} \rangle$ is the linearized loss of the unary loss function $\widetilde{f}_t(\mathbf{x}_t)$ over which each base-learner opti-

**Algorithm 2** Meta OFW Algorithm (**Meta-OFW**).

**Input:** Time horizon $T$; number of base-learners $N$ per eq. (11); step-size pool $\mathcal{H}$ per eq. (12); initial weight of base-learners $\boldsymbol{p}_1$ per eq. (13); learning rate $\epsilon$ for meta-algorithm per eq. (14).
**Output:** Decision $\mathbf{x}_t$ at each time step $t = 1, \ldots, T$.

1: Set $\mathbf{x}_\tau = \mathbf{0}, \forall \tau \leq 0$;
2: Initialize $\mathbf{x}_{1,i} \in \mathcal{X}, \forall i \in \{1, \ldots, N\}$;
3: **for** each time step $t = 1, \ldots, T$ **do**
4:     Receive $\mathbf{x}_{t,i}$ from base-learner $\mathcal{B}_i$ for all $i$;
5:     Output the Decision $\mathbf{x}_t = \sum_{i=1}^{N} p_{t,i}\mathbf{x}_{t,i}$;
6:     Suffer loss $f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t)$;
7:     Observe the loss function $f_t : \mathcal{X}^{m+1} \mapsto \mathbb{R}$;
8:     Construct linearized loss
$$g_t(\mathbf{x}) = \left\langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{x} \right\rangle;$$
9:     Construct the switching-cost-regularized surrogate loss $\ell_t \in \mathbb{R}^N$ with
$$\ell_{t,i} = g_t(\mathbf{x}_{t,i}) + \lambda \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2;$$
10:     Update the weight of base-learners $\boldsymbol{p}_{t+1} \in \Delta_N$ by
$$p_{t+1,i} = \frac{p_{t,i} e^{-\epsilon \ell_{t,i}}}{\sum_{j=1}^{N} p_{t,j} e^{-\epsilon \ell_{t,j}}};$$
11:     Base-learner $\mathcal{B}_i$ updates $\mathbf{x}_{t+1,i}$ with step size $\eta_i$ for all $i$, per **OFW** in Algorithm 1;
12: **end for**

mizes via the **OFW** algorithm;

- $\ell_{t,i} \triangleq g_t(\mathbf{x}_{t,i}) + \lambda \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2$ is a surrogate loss associated with the $i$-th base-learner $\mathcal{B}_i$ —the meta-learner collects $\ell_{t,i}$ for all base-learners, *i.e.*, for all $i \in \{1, \ldots, N\}$, and optimizes $\mathbf{x}_t$ via **Hedge**;
- $p_{t,i}$ is the assigned weight to the $i$-th base-learner $\mathcal{B}_i$ by **Hedge** —each $p_{t,i}$, $i \in \{1, \ldots, N\}$, is used to output **Meta-OFW**'s final decision $\mathbf{x}_t$ as the weighted sum of base-learners' decisions $\mathbf{x}_{t,i}$; *i.e.*, $\mathbf{x}_t = \sum_{i=1}^{N} p_{t,i}\mathbf{x}_{t,i}$;
- $\ell_t \in \mathbb{R}^N$ is the vector whose $i$-th entry is $\ell_{t,i}$;
- $\boldsymbol{p}_t$ is the vector with $i$-th entry as $p_{t,i}$;
- $\alpha \triangleq 2(a+c)$ is a constant introduced for notational simplicity ($a$ and $c$ are per Assumption 3).

**Remark 2** (Unknown Loss Variation $V_T$ Requires Multiple **OFW** Base-Learners). *The multiple* **OFW** *base-learners aim to overcome the challenge of the a priori unknown loss variation $V_T$. To illustrate this, we first consider that $V_T$ is known a priori, and show that a single* **OFW** *suffices to achieve bounded dynamic regret for OCO-M. Then, we consider that $V_T$ is unknown a priori, and show how multiple base-learners with appropriate step sizes $\eta$ can approximate the case where $V_T$ is known a priori. To these ends, we leverage the following dynamic regret bound for OCO-*

*M (Anava et al., 2015, Proof of Theorem 3.1):*

$$\text{Regret}_T^D \leq \underbrace{\sum_{t=1}^{T} \widetilde{f}_t(\mathbf{x}_t) - \sum_{t=1}^{T} \widetilde{f}_t(\mathbf{v}_t)}_{\text{unary cost}}$$

$$+ \lambda \underbrace{\sum_{t=2}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2}_{\text{switching cost}} + \lambda \underbrace{\sum_{t=2}^{T} \|\mathbf{v}_t - \mathbf{v}_{t-1}\|_2}_{\text{path length}},$$

(9)

*which we can simplify to*

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T(1 + V_T + D_T + C_T)}\right), \quad (10)$$

*when $V_T$ is known a priori. Particularly, assume that $\mathbf{x}_t$ is updated by an **OFW** algorithm applied to $\widetilde{f}_1, \ldots, \widetilde{f}_T$ with the $V_T$-dependent step size $\eta_* = \mathcal{O}\left(\sqrt{(1 + V_T)/T}\right)$. Then, eq. (10) results from eq. (9) since the three terms in eq. (9) can be bounded respectively as follows: (i) the unary cost can be bounded by eq. (7) where $\eta = \eta_*$; (ii) the switching cost can be bounded by $\eta_* T D$ due to **OFW**'s line 6 and due to Assumption 1; and (iii) the path length is by definition equal to $C_T$. Then, an application of the Cauchy-Schwarz inequality completes the proof of eq. (10). All in all, when $V_T$ is known a priori, a single **OFW** suffices to achieve bounded dynamic regret for OCO-M.*

*But $V_T$ is unknown a priori since it depends on the loss functions, which are unknown a priori. Instead, an upper bound to $V_T$ is known, specifically, it holds true that $V_T \leq Tc$ under Assumption 3. Leveraging this, we can approximate the case where $V_T$ is known a priori by employing an appropriate number of **OFW** base-learners, each with a different step size, per eq. (11) and eq. (12) below. Intuitively, we can guarantee that way that there exists a base-learner $i$ with step size $\eta_i$ close to the unknown step size $\eta_*$ (the full justification of eq. (11) and eq. (12) is given in Theorem 2's proof in Appendix C.2). The challenge now is to fuse the decisions of the multiple **OFW** base-learners to a final decision $\mathbf{x}_t$.*

**Remark 3** (The Multiple **OFW** Require a **Hedge** Meta-Learner). *The **Hedge** meta-learner in **Meta-OFW** aims to fuse the decisions of the multiple **OFW** base-learners to a final decision $\mathbf{x}_t$. Specifically, the **OFW** base-learners provide multiple decisions at each iteration, the $\mathbf{x}_{t,i}$, $i \in \{1, \ldots, N\}$ (line 4 in Algorithm 2). Then, **Meta-OFW** utilizes the **Hedge** steps in lines 5, 9, and 10 to fuse those decisions to a single decision, aiming to "track" the best base-learner $\mathcal{B}_i$.*

We next formally describe **Meta-OFW**. First, the algorithm specifies the number of base learners, their corresponding

step sizes, and their initial weights as follows, respectively:

$$N = \left\lceil \frac{1}{2} \log_2(1 + \frac{Tc}{\alpha}) \right\rceil + 1 = \mathcal{O}(\log T), \quad (11)$$

$$\mathcal{H} = \left\{ \eta_i \mid \eta_i = 2^{i-1}\sqrt{\frac{\alpha}{\lambda TD}} \leq 1, i \in \{1, \ldots, N\} \right\}, \quad (12)$$

$$p_{1,i} = \frac{1}{i(i+1)} \cdot \frac{N+1}{N}, \text{ for any } i \in \{1, \ldots, N\}. \quad (13)$$

Also, **Meta-OFW** sets the meta-learner's learning rate per the following eq. (14):

$$\epsilon = \sqrt{2/((2\lambda + G)(\lambda + G)D^2T)}, \quad (14)$$

where the dependence on $T$ can be removed by a "doubling trick" (Cesa-Bianchi et al., 1997), similarly to how **Meta-OFW** copes with the unknown $V_T$,

At each iteration $t = 1, \ldots, T$, **Meta-OFW** receives the intermediate decisions $\mathbf{x}_{t,i}$ from all the base-learners $\mathcal{B}_i$, $i \in \{1, \ldots, N\}$ (line 4) to fuse them into a final decision $\mathbf{x}_t = \sum_{i=1}^{N} p_{t,i} \mathbf{x}_{t,i}$ (line 5). Then, **Meta-OFW** suffers a loss of $f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t)$ (lines 6-7). Afterwards, **Meta-OFW** constructs the linearized loss $g_t(\mathbf{x})$ and switching-cost-regularized loss $\ell_t$ (lines 8-9). To this end, **Meta-OFW** needs to evaluate only once the gradient $\nabla \widetilde{f}_t(\mathbf{x}_t)$. Finally, the meta-learner and base-learners update the weights $p_{t+1}$ and $\mathbf{x}_{t+1,i}$ for the next iteration (lines 10-11).

## 5. Dynamic Regret Guarantees of Meta-OFW

To present **Meta-OFW**'s dynamic regret bound, we define:

- $D_{T,i} \triangleq \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_{t,i}) - \nabla f_{t-1}(\mathbf{x}_{t-1,i})\|_2^2$ is the gradient variation associated with the base-learner $i$;
- $\bar{D}_T \triangleq \max_{i \in \{1, \ldots, N\}} D_{T,i}$ is the upper bound for $D_{T,i}$.

We present **Meta-OFW**'s dynamic regret bound against any comparator sequence (Theorem 2). Particularly, the bound below holds true, even if the loss variation $V_T$, gradient variation $\bar{D}_T$, and path length $C_T$ are unknown to **Meta-OFW**.

**Theorem 2** (Dynamic Regret Bound of **Meta-OFW**). *For any comparator sequence $(\mathbf{v}_1, \ldots, \mathbf{v}_T) \in \mathcal{X}^T$, **Meta-OFW** achieves a dynamic regret $\text{Regret}_T^D$ that enjoys the bound:*

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T(1 + V_T + \bar{D}_T + C_T)}\right). \quad (15)$$

The dependency on $V_T$ and $\bar{D}_T$ results from **OFW** being a base-learner in **Meta-OFW**; similar dependencies, due to projection-free subroutines in online algorithms, have been observed in the literature: see, *e.g.*, (Kalhan et al., 2021) and the references in Table 1. The dependency on $\bar{D}_T$,

instead of $D_T$ in Theorem 1, is to upper bound the gradient variation $D_{T,i}$ such that the base-learner $i$ with step size close to the unknown step size $\eta^\star$ (Remark 2) satisfies $D_{T,i} \leq \bar{D}_T$.

The dependency on $C_T$ results from the sequence of comparators being time-varying. Specifically, (Zhang et al., 2018) proved that any optimal dynamic regret bound for OCO is $\Omega\left(\sqrt{T(1+C_T)}\right)$, and thus the bound necessarily depends on $C_T$ in the worst case.

**Remark 4** (Trade-Off of Projection-Free Efficiency with Regret Optimality). *(Zhao et al., 2022) prove that the optimal dynamic regret for OCO-M is $\Omega(\sqrt{T(1+C_T)})$, and provide a projection-based algorithm using **OGD** that matches this bound. In contrast, **Meta-OFW**'s regret bound in Theorem 2 cannot match the bound $\Omega(\sqrt{T(1+C_T)})$ due to the presence of $V_T$ and $\bar{D}_T$ in eq. (15). But **Meta-OFW** is projection-free and thus is more efficient than the **OGD**-based algorithm in (Zhao et al., 2022) (Hazan & Kale, 2012). All in all, the dependence of eq. (15) on $V_T$ and $\bar{D}_T$ is the regret suboptimality cost we pay in this paper to solve OCO-M efficiently via the projection-free **OFW**.*

# 6. Application to Non-Stochastic Control

We apply **Meta-OFW** to the online non-stochastic control problem (Agarwal et al., 2019), and present a projection-free controller with memory (Algorithm 3), and with bounded dynamic regret against any linear time-varying feedback control policy (Theorem 3). The results of the numerical evaluations are present in Appendix E.

## 6.1. The Non-Stochastic Control Problem

We consider Linear Time-Varying systems of the form

$$x_{t+1} = A_t x_t + B_t u_t + w_t, \quad t = 0, \ldots, T, \quad (16)$$

where $x_t \in \mathbb{R}^{d_x}$ is the state of the system, $u_t \in \mathbb{R}^{d_u}$ is the control input, and $w_t \in \mathbb{R}^{d_x}$ is the process noise. The system and input matrices, $A_t$ and $B_t$, respectively, are known.

At each time step $t$, the controller chooses a control action $u_t$ and then suffers a loss $c_t(x_t, u_t)$. The loss function $c_t$ is revealed to the controller only after the controller has chosen the control action $u_t$, similarly to the OCO setting.

**Assumption 6** (Convex and Bounded Loss Function with Bounded Gradient). *The cost function $c_t(x_t, u_t) : \mathbb{R}^{d_x} \times \mathbb{R}^{d_u} \mapsto \mathbb{R}$ is convex in $x_t$ and $u_t$. Further, when $\|x\|_2 \leq D$, $\|u\|_2 \leq D$ for some $D > 0$, then $|c_t(x,u)| \leq \beta D^2$ and $\|\nabla_x c_t(x,u)\|_2 \leq G_c D, \|\nabla_u c_t(x,u)\|_2 \leq G_c D$, for given positive numbers $\beta$ and $G_c$.*

**Assumption 7** (Bounded System Matrices and Noise). *The system matrices and noise are bounded, i.e., $\|A_t\|_{\text{op}} \leq \kappa_A$,*

$\|B_t\|_{\text{op}} \leq \kappa_B$, *and $\|w_t\|_2 \leq W$ for given positive numbers $\kappa_A$, $\kappa_B$, and $W$, where $\|\cdot\|_{\text{op}}$ is the operator norm.*

Per Assumption 7, we assume no stochastic model for the process noise $w_t$: the noise may even be adversarial, subject to the bounds prescribed by $W$.

**Problem 2** (Non-Stochastic Control (NSC) Problem). *At each time step $t = 0, \ldots, T$, first a control action $u_t$ is chosen; then, a loss function $c_t : \mathbb{R}^{d_x} \times \mathbb{R}^{d_u} \mapsto \mathbb{R}$ is revealed and the system suffers a loss $c_t(x_t, u_t)$. The goal is to minimize the dynamic policy regret defined below.*

**Definition 2** (Dynamic Policy Regret). *We define the dynamic policy regret as*

$$\text{Regret-NSC}_T^D = \sum_{t=0}^{T} c_t(x_t, u_t) - \sum_{t=0}^{T} c_t(x_t^*, u_t^*), \quad (17)$$

*where (i) both sums in eq. (17) are evaluated with the same noise $\{w_0, \ldots, w_T\}$, which is the noise experienced by the system during its evolution per the control input $\{u_0, \ldots, u_T\}$, (ii) $u_t^* = -K_t^* x_t^*$ is the optimal linear feedback control input in hindsight, i.e., the optimal input given a priori knowledge of $c_t$ and of the realized noise $w_t$, and (iii) $x_t^*$ is the state reached by applying the sequence of optimal control inputs $\{u_0^*, \ldots, u_{t-1}^*\}$.*

**Reduction to OCO-M.** We present the reduction of the non-stochastic control problem to OCO-M, following (Agarwal et al., 2019).

Per eq. (16), $x_t$ depends on the control actions chosen in the past, *i.e.*, $\{u_0, \ldots, u_{t-1}\}$, and similarly, the control action $u_t$ depends on $x_{t-1}$, *i.e.*, $\{u_0, \ldots, u_{t-2}\}$. To reduce the non-stochastic control problem to OCO-M, there are thus 2 challenges: (i) we need a control parameterization such that the cost function $c_t(x_t, u_t)$ is convex in the parameters of the control actions $\{u_0, \ldots, u_{t-1}\}$, since $c_t(x_t, u_t)$ is implicitly a function of $\{u_0, \ldots, u_{t-1}\}$ via $u_t$; and, similarly, (ii) we need the memory length of $c_t(x_t, u_t)$, *i.e.*, its implicit dependence on the past control inputs $\{u_0, \ldots, u_{t-1}\}$, to stop growing as $t$ increases; that is, we need $c_t(x_t, u_t)$ to instead depend on the most recent control inputs only, in particular, on $\{u_{t-m}, \ldots, u_t\}$ for memory length $m$. To address these challenges, (Agarwal et al., 2019) propose the *Disturbance-Action Control* policy and the notion of *truncated loss*.

**Definition 3** (Disturbance-Action Control Policy). *A Disturbance-Action Control (DAC) policy $\pi_t(K_t, M_t)$ chooses the control action $u_t$ at state $x_t$ as $u_t = -K_t x_t + \sum_{i=1}^{H} M_t^{[i-1]} w_{t-i}$,[5] where $M_t = (M_t^{[0]}, \ldots, M_t^{[H-1]})$ with $\left\|M_t^{[i]}\right\|_{\text{op}} \leq \kappa_B \kappa^3 (1-\gamma)^i$ and horizon $H \geq 1$, $K_t$ is*

---

[5]The DAC policy depends on the past noise, which can be obtained from eq. (16) once the next state is observed; specifically, at time $t+1$, it holds true that $w_t = x_{t+1} - A_t x_t - B_t u_t$.

**Algorithm 3 Meta-OFW** for Non-Stochastic Control.

---

**Input:** Time horizon $T$; number of base-learners $N$ per eq. (18); step size pool $\mathcal{H}$ per eq. (19); initial weight of base-learners $\boldsymbol{p}_0$ per eq. (20); learning rate $\epsilon$ of meta-algorithm per eq. (21).

**Output:** Control $u_t$ at each time step $t = 1, \ldots, T$.

1: Set $M_\tau = \mathbf{0}$ and $w_\tau = 0, \forall \tau < 0$;
2: Initialize $M_{0,i} \in \mathcal{M}, \forall i \in \{1, \ldots, N\}$;
3: **for** each time step $t = 0, \ldots, T$ **do**
4:    Receive $M_{t,i}$ from base-learner $\mathcal{B}_i$ for all $i$;
5:    Calculate $M_t = \sum_{i=1}^N p_{t,i} M_{t,i}$;
6:    Output $u_t = -K_t x_t + \sum_{i=1}^H M_t^{[i-1]} w_{t-i}$;
7:    Observe the loss function $c_t : \mathbb{R}^{d_x} \times \mathbb{R}^{d_u} \mapsto \mathbb{R}$ and suffer the loss $c_t(x_t, u_t)$;
8:    Construct the truncated loss $f_t(M_{t-H-1}, \ldots, M_t)$ : $\mathcal{M}^{H+2} \mapsto \mathbb{R}$ ;
9:    Construct the linearized loss
$$g_t(M) = \left\langle \nabla_M \widetilde{f}_t(M_t), M \right\rangle_{\mathrm{F}} ;$$
10:   Construct the switching-cost-regularized surrogate loss $\boldsymbol{\ell}_t \in \mathbb{R}^N$ with
$$\ell_{t,i} = g_t(M_{t,i}) + \zeta \|M_{t,i} - M_{t-1,i}\|_{\mathrm{F}} ;$$
11:   Update the weight of base-learners $\boldsymbol{p}_{t+1} \in \Delta_N$ by
$$p_{t+1,i} = \frac{p_{t,i} e^{-\epsilon \ell_{t,i}}}{\sum_{j=1}^N p_{t,j} e^{-\epsilon \ell_{t,j}}} ;$$
12:   **for** each base-learner $\mathcal{B}_i$ **do**
13:      Compute
$$M'_{t,i} = \arg \min_{M \in \mathcal{M}} \left\langle \nabla_M \widetilde{f}_t(M_t), M \right\rangle_{\mathrm{F}} ;$$
14:      Update $M_{t+1,i} = (1 - \eta_i) M_{t,i} + \eta_i M'_{t,i}$;
15:   **end for**
16:   Observe the state $x_{t+1}$ and calculate the noise $w_t = x_{t+1} - A_t x_t - B_t u_t$;
17: **end for**

---

$a (\kappa, \gamma)-$strongly stable matrix which is calculated given $A_t$ and $B_t$, and $w_\tau = 0$ for all $\tau < 0$.

Per Proposition 1 in Appendix D.4 (Gradu et al., 2020), $x_t$ and $u_t$ are linear in $\{M_0, \ldots, M_t\}$; therefore, the cost function $c_t(x_t, u_t)$ is convex in $\{M_0, \ldots, M_t\}$.

To present the notion of *truncated loss*, we use the notation:

- $x_t (M_{0:t-1})$ is the state reached by applying the DAC policy $\{\pi_\tau (K_\tau, M_\tau)\}_{\tau=0,\ldots,t-1}$;
- $u_t (M_{0:t})$ is the control action at state $x_t (M_{0:t-1})$ per the DAC policy $\pi_t(K_t, M_t)$;
- $y_t (M_{t-1-H:t-1})$ is the state reached from $x_{t-1-H} = 0$ by applying $\{\pi_\tau (K_\tau, M_\tau)\}_{\tau=t-H-1,\ldots,t-1}$ and experiencing the noise sequence $\{w_\tau\}_{\tau=t-H-1,\ldots,t-1}$;
- $v_t (M_{t-1-H:t})$ is the control input that would have been executed if the state at time $t$ was the $y_t (M_{t-1-H:t-1})$.

**Definition 4** (Truncated Loss). *Given DAC policies* $\{\pi_\tau (K_\tau, M_\tau)\}_{\tau=0,\ldots,t}$ *with memory length* $H$, *the induced truncated loss* $f_t : \mathcal{M}^{H+2} \mapsto \mathbb{R}$ *is defined as*

$$f_t (M_{t-1-H:t}) \triangleq c_t (y_t (M_{t-1-H:t-1}), v_t (M_{t-1-H:t})) .$$

Thereby, the truncated loss $f_t (M_{t-1-H:t})$ depends only on the last $H + 2$ time steps of the DAC policy. That is, $f_t$ has a fixed memory length $H + 2$, for all $t = 1, \ldots, T$.

All in all, Problem 2 can be reduced to OCO-M when the decision variables are the $M_t$, and the loss functions are the truncated losses $f_t (M_{t-1-H:t})$, for all $t = 1, \ldots, T$.

### 6.2. Meta-OFW for Online Non-Stochastic Control

We present **Meta-OFW**'s application to the online non-stochastic control problem (Algorithm 3). Particularly, Algorithm 3 initializes the number of base-learners, their corresponding step sizes, and their initial weights, per the following equations, similarly to **Meta-OFW**:

$$N = \left\lceil \frac{1}{2} \log_2 (\frac{2\beta D^2 T + \phi}{\sigma}) \right\rceil + 1 = \mathcal{O}(\log T), \quad (18)$$

$$\mathcal{H} = \left\{ \eta_i \mid \eta_i = 2^{i-1} \sqrt{\frac{\sigma}{\zeta T D_f}} \leq 1, i \in \{1, \ldots, N\} \right\}, \quad (19)$$

$$p_{0,i} = \frac{1}{i(i+1)} \cdot \frac{N+1}{N}, \text{ for any } i \in \{1, \ldots, N\}, \quad (20)$$

where $\sigma \triangleq 4\beta D^2$, $\phi \triangleq \sigma + 2\beta D^2$, $\zeta \triangleq (H+2)^2 L_f$, and $L_f, G_f$ defined as in Lemma 9 in Appendix D.7.

The algorithm also sets the step size of meta-learner as

$$\epsilon = \sqrt{2/\left((2\zeta + G_f)(\zeta + G_f)D_f^2 T\right)}. \quad (21)$$

At each iteration $t$, Algorithm 3 receives $M_{t,i}$ from all base-learners (line 4). Then, Algorithm 3 calculates $M_t = \sum_{i=1}^N p_{t,i} M_{t,i}$ and outputs the control actions $u_t = -K_t x_t + \sum_{i=1}^H M_t^{[i-1]} w_{t-i}$ (lines 5-6), after which the cost function is revealed and the algorithm suffers a loss of $c_t(x_t, u_t)$ (line 7). Next, Algorithm 3 constructs the truncated loss $f_t(M_{t-H-1}, \ldots, M_t)$, linearized loss $g_t(M)$, and switching-cost-regularized loss $\boldsymbol{\ell}_t$ (lines 8-10). The meta-learner and base-learners update the weights $\boldsymbol{p}_{t+1}$ and $M_{t+1,i}$ for the next iteration (lines 11-15). Finally, the noise $w_t$ is calculated when $x_{t+1}$ is observed (line 16).

### 6.3. Dynamic Regret Guarantee of Algorithm 3

**Theorem 3** (Dynamic Policy Regret Bound of Algorithm 3). *Algorithm 3 ensures that* [6]

$$\text{Regret-NSC}_T^D \leq \tilde{\mathcal{O}}\left(\sqrt{T\left(1 + V_T + \bar{D}_T + C_T\right)}\right). \tag{22}$$

**Remark 5** (Novelty of Theorem 3). *Theorem 3 guarantees a dynamic regret bound against an optimal time-varying linear feedback policy in hindsight, i.e., against $\{\pi_\tau(K_\tau^*, 0)\}_{\tau=0,\dots,t}$, per Definition 2. This is different than competing against an optimal time-varying DAC policy $\{\pi_\tau(K_\tau, M_\tau^*)\}_{\tau=0,\dots,t}$ with pre-specified stabilizing control gains $K_\tau$ as in (Zhao et al., 2022), or an optimal time-invariant linear feedback policy $\{\pi(K^*, 0)\}$ over the entire horizon or any time interval as in (Agarwal et al., 2019; Li et al., 2021; Gradu et al., 2020; Simchowitz et al., 2020; Zhang et al., 2022). To achieve this, we show a DAC policy $\{\pi_\tau(K_\tau, M_\tau)\}_{\tau=0,\dots,t}$ can approximate any time-varying linear feedback policy (Proposition 2 in Appendix D.5) without the need to specify to the optimal policy any stabilizing feedback control gains.*

## 7. Conclusion

We provided **Meta-OFW** (Algorithm 2), the first projection-free algorithm with bounded dynamic regret for OCO with memory in time-varying environments (Theorem 2). To develop **Meta-OFW**, we employed the projection-free algorithm **OFW** along with **Hedge**. Further, we applied **Meta-OFW** to the online non-stochastic control problem to control linear time-varying systems that are corrupted with unknown and unpredictable noise (Algorithm 3). We thus developed a projection-free controller with memory and bounded dynamic regret against any linear time-varying control policy (Theorem 3), instead of against only static linear control policies. To this end, we also proved that the DAC policy class (Agarwal et al., 2019) can approximate linear time-varying feedback controllers (Proposition 2 in Appendix D.5).

## References

Agarwal, N., Bullins, B., Hazan, E., Kakade, S., and Singh, K. Online control with adversarial disturbances. In *International Conference on Machine Learning (ICML)*, pp. 111–119, 2019.

Anava, O., Hazan, E., and Mannor, S. Online learning for adversaries with memory: price of past mistakes. *Advances in Neural Information Processing Systems (NeurIPS)*, 28, 2015.

Åström, K. J. *Introduction to stochastic control theory*. Courier Corporation, 2012.

Besbes, O., Gur, Y., and Zeevi, A. Non-stationary stochastic optimization. *Operations Research*, 63(5):1227–1244, 2015.

Cesa-Bianchi, N. and Lugosi, G. *Prediction, learning, and games*. Cambridge university press, 2006.

Cesa-Bianchi, N., Freund, Y., Haussler, D., Helmbold, D. P., Schapire, R. E., and Warmuth, M. K. How to use expert advice. *Journal of the ACM (JACM)*, 44(3): 427–485, 1997.

Chang, T.-J. and Shahrampour, S. Unconstrained online optimization: Dynamic regret analysis of strongly convex and smooth problems. *arXiv preprint:2006.03912*, 2020.

Chen, G. and Teboulle, M. Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993.

Duchi, J. C., Shalev-Shwartz, S., Singer, Y., and Tewari, A. Composite objective mirror descent. In *Conference on Learning Theory (COLT)*, volume 10, pp. 14–26. PMLR, 2010.

Frank, M. and Wolfe, P. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3(1-2): 95–110, 1956.

Freund, Y. and Schapire, R. E. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1): 119–139, 1997.

Garber, D. and Hazan, E. Faster rates for the frank-wolfe method over strongly-convex sets. In *International Conference on Machine Learning*, pp. 541–549. PMLR, 2015.

Garber, D. and Kretzu, B. New projection-free algorithms for online convex optimization with adaptive regret guarantees. *arXiv preprint arXiv:2202.04721*, 2022.

Gradu, P., Hazan, E., and Minasyan, E. Adaptive regret for control of time-varying dynamics. *arXiv preprint:2007.04393*, 2020.

Hazan, E. and Kale, S. Projection-free online learning. *arXiv preprint:1206.4657*, 2012.

Hazan, E. and Seshadhri, C. Adaptive algorithms for online decision problems. In *Electronic colloquium on computational complexity (ECCC)*, volume 14, 2007.

Hazan, E., Kakade, S., and Singh, K. The nonstochastic control problem. In *Algorithmic Learning Theory (ALT)*, pp. 408–421, 2020.

---

[6] The path length is defined as $C_T \triangleq \sum_{t=2}^T \|M_{t-1}^* - M_t^*\|_F$.

Hazan, E. et al. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3-4): 157–325, 2016.

Jadbabaie, A., Rakhlin, A., Shahrampour, S., and Sridharan, K. Online optimization: Competing with dynamic comparators. In *Artificial Intelligence and Statistics (AISTATS)*, pp. 398–406. PMLR, 2015.

Jaggi, M. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International Conference on Machine Learning*, pp. 427–435. PMLR, 2013.

Kalhan, D. S., Bedi, A. S., Koppel, A., Rajawat, K., Hassani, H., Gupta, A. K., and Banerjee, A. Dynamic online learning via frank-wolfe algorithm. *IEEE Transactions on Signal Processing (TSP)*, 69:932–947, 2021.

Kretzu, B. and Garber, D. Revisiting projection-free online learning: the strongly convex case. In *International Conference on Artificial Intelligence and Statistics*, pp. 3592–3600. PMLR, 2021.

Li, Y., Das, S., and Li, N. Online optimal control with affine constraints. In *AAAI Conference on Artificial Intelligence (AAAI)*, volume 35, pp. 8527–8537, 2021.

Mokhtari, A., Shahrampour, S., Jadbabaie, A., and Ribeiro, A. Online optimization in dynamic environments: Improved regret rates for strongly convex problems. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pp. 7195–7201, 2016.

Rockafellar, R. T. Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization*, 14(5):877–898, 1976.

Romero, A., Penicka, R., and Scaramuzza, D. Time-optimal online replanning for agile quadrotor flight. *IEEE Robotics and Automation Letters*, 7(3):7730–7737, 2022.

Shalev-Shwartz, S. et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.

Shi, G., Shi, X., O'Connell, M., Yu, R., Azizzadenesheli, K., Anandkumar, A., Yue, Y., and Chung, S.-J. Neural lander: Stable drone landing control using learned dynamics. In *International Conference on Robotics and Automation (ICRA)*, pp. 9784–9790, 2019.

Simchowitz, M., Singh, K., and Hazan, E. Improper learning for non-stochastic control. In *Conference on Learning Theory (COLT)*, pp. 3320–3436, 2020.

Vishnoi, N. K. *Algorithms for convex optimization*. Cambridge University Press, 2021.

Wan, Y. and Zhang, L. Projection-free online learning over strongly convex sets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 10076–10084, 2021.

Wan, Y., Zhang, L., and Song, M. Improved dynamic regret for online frank-wolfe. *arXiv preprint arXiv:2302.05620*, 2023.

Weinberger, M. J. and Ordentlich, E. On delayed prediction of individual sequences. *IEEE Transactions on Information Theory (TIT)*, 48(7):1959–1976, 2002.

Xu, Z., Zhou, H., and Tzoumas, V. Online submodular coordination with bounded tracking regret: Theory, algorithm, and applications to multi-robot coordination. *IEEE Robotics and Automation Letters (RAL)*, 2023.

Yang, T., Zhang, L., Jin, R., and Yi, J. Tracking slowly moving clairvoyant: Optimal dynamic regret of online learning with true and noisy gradient. In *International Conference on Machine Learning (ICML)*, pp. 449–457. PMLR, 2016.

Zhang, L. Online learning in changing environments. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence (IJCAI)*, pp. 5178–5182, 2020.

Zhang, L., Lu, S., and Zhou, Z.-H. Adaptive online learning in dynamic environments. *Advances in Neural Information Processing Systems (NeurIPS)*, 31, 2018.

Zhang, Z., Cutkosky, A., and Paschalidis, I. Adversarial tracking control via strongly adaptive online learning with memory. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 8458–8492. PMLR, 2022.

Zhao, P., Yan, Y.-H., Wang, Y.-X., and Zhou, Z.-H. Non-stationary online learning with memory and non-stochastic control. *arXiv preprint arXiv:2102.03758*, 2021.

Zhao, P., Wang, Y.-X., and Zhou, Z.-H. Non-stationary online learning with memory and non-stochastic control. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 2101–2133. PMLR, 2022.

Zhao, Y., Zhao, Q., Zhang, X., Zhu, E., Liu, X., and Yin, J. Understand dynamic regret with switching cost for online decision making. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(3):1–21, 2020.

Zinkevich, M. Online convex programming and generalized infinitesimal gradient ascent. In *Interna. Conf. on Machine Learning (ICML)*, pp. 928–936, 2003.

# A. Dynamic Regret Analysis of OFW in OCO without Memory

**Theorem 4** (Dynamic Regret Bound of **OFW**). *Consider the OCO problem with no memory, i.e., Problem 1 with $m = 0$. Under Assumption 1 and Assumption 2,* **OFW** *achieves against any sequence of comparators $(\mathbf{v}_1, \ldots, \mathbf{v}_T) \in \mathcal{X}^T$ the dynamic regret*

$$\text{Regret}_T^D \leq \mathcal{O}\left(\frac{1 + V_T}{\eta} + \sqrt{T D_T}\right). \tag{23}$$

*Under step size $\eta = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$, we have the following dynamic regret bound,*

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T}\left(1 + V_T + \sqrt{D_T}\right)\right). \tag{24}$$

*Further, select $\eta = \sqrt{\frac{c}{b}}$ with $c < b$, we have the following dynamic regret bound,*

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T(V_T + D_T)}\right). \tag{25}$$

## A.1. Proof of Theorem 4

*Proof.* The proof follows the steps of (Kalhan et al., 2021, Proof of Theorem 1) but removes the assumption that the loss function $f_t$ must be smooth per the original proof in (Kalhan et al., 2021).

Additionally, the proof concludes with the novel bound in eq. (25), which is enabled with a constant step size (eq. (36)) under Assumption 3.

In more detail, to prove Theorem 1, we take summation on both sides from $t = 1$ to $T$ of eq. (37) in Lemma 1:

$$\sum_{t=1}^{T}(f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t)) \leq \sum_{t=1}^{T}\left(f_{t,t-1}^{\text{sup}} + f_{t-1}(\mathbf{v}_{t-1}) - f_t(\mathbf{v}_t)\right)$$
$$+ (1 - \eta)\left(\sum_{t=1}^{T-1}(f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t)) + f_0(\mathbf{x}_0) - f_0(\mathbf{v}_0)\right) \tag{26}$$
$$+ \eta D \sum_{t=1}^{T}\|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2,$$

where $f_{t,t-1}^{\text{sup}}$ is defined in Lemma 1. Moving $(1 - \eta)\sum_{t=1}^{T-1}(f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t))$ to the left side of eq. (26), and noting $\sum_{t=1}^{T}(f_{t-1}(\mathbf{v}_{t-1}) - f_t(\mathbf{v}_t)) = f_0(\mathbf{v}_0) - f_T(\mathbf{v}_T)$, we get

$$\eta \sum_{t=1}^{T-1}(f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t)) + (f_T(\mathbf{x}_T) - f_T(\mathbf{v}_T)) \leq \sum_{t=1}^{T} f_{t,t-1}^{\text{sup}} + (1 - \eta)(f_0(\mathbf{x}_0) - f_0(\mathbf{v}_0)) + f_0(\mathbf{v}_0) - f_T(\mathbf{v}_T)$$
$$+ \eta D \sum_{t=1}^{T}\|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2. \tag{27}$$

Subtracting $(1 - \eta)(f_T(\mathbf{x}_T) - f_T(\mathbf{v}_T))$ from both side of eq. (27), we obtain

$$\eta \text{Regret}_T^D \leq \sum_{t=1}^{T} f_{t,t-1}^{\text{sup}} + (1 - \eta)(f_0(\mathbf{x}_0) - f_T(\mathbf{x}_T) - f_0(\mathbf{v}_0) + f_T(\mathbf{v}_T))$$
$$+ f_0(\mathbf{v}_0) - f_T(\mathbf{v}_T) + \eta D \sum_{t=1}^{T}\|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2$$
$$= \sum_{t=1}^{T} f_{t,t-1}^{\text{sup}} + \eta(-f_0(\mathbf{x}_0) + f_T(\mathbf{x}_T) + f_0(\mathbf{v}_0) - f_T(\mathbf{v}_T)) \tag{28}$$
$$+ f_0(\mathbf{x}_0) - f_T(\mathbf{x}_T) + \eta D \sum_{t=1}^{T}\|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2.$$

By Assumption 3 and the Cauchy-Schwarz inequality, it holds true that: $f_0(\mathbf{x}_0) - f_T(\mathbf{x}_T) \leq 2(a+c)$, and $-f_0(\mathbf{x}_0) + f_T(\mathbf{x}_T) + f_0(\mathbf{v}_0) - f_T(\mathbf{v}_T) \leq 4(a+c)$. Divide both sides of eq. (28) by $\eta$ and substitute the following inequality, which holds true due to the Cauchy-Schwarz inequality,

$$\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2 \leq \sqrt{T \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2} = \sqrt{TD_T} \tag{29}$$

into eq. (28) with the definition of $V_T$ to obtain

$$\text{Regret}_T^D \leq \frac{1}{\eta} V_T + \frac{2}{\eta}(a+c) + D\sqrt{TD_T} + 4(a+c). \tag{30}$$

Selecting $\eta = \mathcal{O}(\frac{1}{\sqrt{T}})$ we get

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T}\left(1 + V_T + \sqrt{D_T}\right)\right). \tag{31}$$

Further, by selecting $\eta = \sqrt{\frac{c}{b}}$ with $c < b$, the dynamic regret in eq. (30) becomes

$$
\begin{aligned}
\text{Regret}_T^D &\leq \sqrt{\frac{b}{c}} V_T + D\sqrt{TD_T} + \left(\sqrt{\frac{4b}{c}} + 4\right)(a+c) \\
&= \sqrt{Tb}\frac{V_T}{\sqrt{Tc}} + D\sqrt{TD_T} + \left(\sqrt{\frac{4b}{c}} + 4\right)(a+c) \\
&= \sqrt{TV_T b}\frac{\sqrt{V_T}}{\sqrt{Tc}} + D\sqrt{TD_T} + \left(\sqrt{\frac{4b}{c}} + 4\right)(a+c).
\end{aligned}
\tag{32}
$$

From Assumption 3, we have the following bound of $V_T$,

$$0 \leq V_T = \sum_{t=1}^{T} \sup_{\mathbf{x} \in \mathcal{X}} |f_t(\mathbf{x}) - f_{t-1}(\mathbf{x})| \leq Tc, \tag{33}$$

which implies $\frac{\sqrt{V_T}}{\sqrt{Tc}} \leq 1$. Substituting it into eq. (32) gives

$$\text{Regret}_T^D \leq \sqrt{TV_T b} + D\sqrt{TD_T} + \left(\sqrt{\frac{4b}{c}} + 4\right)(a+c). \tag{34}$$

Hence, the dynamic regret is bounded as

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T}\left(\sqrt{V_T} + \sqrt{D_T}\right)\right). \tag{35}$$

Equivalently, due to the Cauchy-Schwarz inequality, the bound can be written as

$$\text{Regret}_T^D \leq \mathcal{O}\left(\sqrt{T(V_T + D_T)}\right), \tag{36}$$

$\square$

## A.2. Proof of Lemma 1

**Lemma 1.** *Under Assumption 1 and Assumption 2, Algorithm 1 satisfies the following descent relations against any sequence of comparators $(\mathbf{v}_1, \ldots, \mathbf{v}_T) \in \mathcal{X}^T$,*

$$
\begin{aligned}
f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t) \leq f_{t,t-1}^{sup} &+ (1-\eta)(f_{t-1}(\mathbf{x}_{t-1}) - f_{t-1}(\mathbf{v}_{t-1})) \\
&+ f_{t-1}(\mathbf{v}_{t-1}) - f_t(\mathbf{v}_t) + \eta D\|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2
\end{aligned}
\tag{37}
$$

*where $f_{t,t-1}^{sup} \triangleq \sup_{\mathbf{x} \in \mathcal{X}} |f_t(\mathbf{x}) - f_{t-1}(\mathbf{x})|$ is the instantaneous maximum cost variation.*

*Proof.* The proof follows similar steps of (Kalhan et al., 2021, Proof of Lemma 1) but removes the assumption that the loss function $f_t$ must be smooth per the original proof in (Kalhan et al., 2021).

By convexity of $f_t(\cdot)$, we have

$$f_t(\mathbf{x}_t) \leq f_t(\mathbf{x}_{t-1}) + \langle \nabla f_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}_{t-1} \rangle. \tag{38}$$

Substituting the update step $\mathbf{x}_t = (1-\eta)\mathbf{x}_{t-1} + \eta\mathbf{x}'_{t-1}$ in Algorithm 1, *i.e.*, $\mathbf{x}_t - \mathbf{x}_{t-1} = \eta(\mathbf{x}'_{t-1} - \mathbf{x}_{t-1})$, into eq. (38),

$$f_t(\mathbf{x}_t) \leq f_t(\mathbf{x}_{t-1}) + \eta\langle \nabla f_t(\mathbf{x}_t), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle. \tag{39}$$

Adding and subtracting the terms $\eta\langle \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle$ to the right hand side of eq. (39), we obtain

$$\begin{aligned} f_t(\mathbf{x}_t) \leq &f_t(\mathbf{x}_{t-1}) + \eta\langle \nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle \\ &+ \eta\langle \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle. \end{aligned} \tag{40}$$

Next, by the optimality condition of $\mathbf{x}'_{t-1}$, *i.e.*,

$$\langle \mathbf{x}'_{t-1}, \nabla f_{t-1}(\mathbf{x}_{t-1}) \rangle = \min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \nabla f_{t-1}(\mathbf{x}_{t-1}) \rangle \leq \langle \mathbf{v}_{t-1}, \nabla f_{t-1}(\mathbf{x}_{t-1}) \rangle, \tag{41}$$

and substituting into eq. (40) leads to

$$\begin{aligned} f_t(\mathbf{x}_t) \leq &f_t(\mathbf{x}_{t-1}) + \eta\langle \nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle \\ &+ \eta\langle \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{v}_{t-1} - \mathbf{x}_{t-1} \rangle. \end{aligned} \tag{42}$$

By convexity of $f_{t-1}(\cdot)$ in eq. (42), we have

$$\begin{aligned} f_t(\mathbf{x}_t) \leq &f_t(\mathbf{x}_{t-1}) + \eta\langle \nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle \\ &+ \eta(f_{t-1}(\mathbf{v}_{t-1}) - f_{t-1}(\mathbf{x}_{t-1})). \end{aligned} \tag{43}$$

Subtracting $f_t(\mathbf{v}_t)$ from both sides of eq. (43) gives

$$\begin{aligned} f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t) \leq &f_t(\mathbf{x}_{t-1}) - f_t(\mathbf{v}_t) + \eta\langle \nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle \\ &+ \eta(f_{t-1}(\mathbf{v}_{t-1}) - f_{t-1}(\mathbf{x}_{t-1})). \end{aligned} \tag{44}$$

Next, consider the term $f_t(\mathbf{x}_{t-1}) - f_t(\mathbf{v}_t)$ from the right hand side of eq. (44). We can bound it as follows:

$$\begin{aligned} f_t(\mathbf{x}_{t-1}) - f_t(\mathbf{v}_t) &= f_t(\mathbf{x}_{t-1}) - f_{t-1}(\mathbf{x}_{t-1}) + f_{t-1}(\mathbf{x}_{t-1}) - f_{t-1}(\mathbf{v}_{t-1}) + f_{t-1}(\mathbf{v}_{t-1}) - f_t(\mathbf{v}_t) \\ &\leq f_{t,t-1}^{\sup} + f_{t-1}(\mathbf{x}_{t-1}) - f_{t-1}(\mathbf{v}_{t-1}) + f_{t-1}(\mathbf{v}_{t-1}) - f_t(\mathbf{v}_t), \end{aligned} \tag{45}$$

where $f_{t,t-1}^{\sup} \triangleq \sup_{\mathbf{x} \in \mathcal{X}}|f_t(\mathbf{x}) - f_{t-1}(\mathbf{x})|$. Substituting eq. (45) into eq. (44), we obtain

$$\begin{aligned} f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t) \leq &f_{t,t-1}^{\sup} + \eta\langle \nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x}'_{t-1} - \mathbf{x}_{t-1} \rangle \\ &+ (1-\eta)(f_{t-1}(\mathbf{x}_{t-1}) - f_{t-1}(\mathbf{v}_{t-1})) + f_{t-1}(\mathbf{v}_{t-1}) - f_t(\mathbf{v}_t). \end{aligned} \tag{46}$$

Applying the Cauchy-Schwarz inequality, we get

$$\begin{aligned} f_t(\mathbf{x}_t) - f_t(\mathbf{v}_t) \leq &f_{t,t-1}^{\sup} + \eta\|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2\|\mathbf{x}'_{t-1} - \mathbf{x}_{t-1}\|_2 \\ &+ (1-\eta)(f_{t-1}(\mathbf{x}_{t-1}) - f_{t-1}(\mathbf{v}_{t-1})) + f_{t-1}(\mathbf{v}_{t-1}) - f_t(\mathbf{v}_t). \end{aligned} \tag{47}$$

Utilizing now Assumption 1 provides the result in Lemma 1. $\qquad\square$

## B. Dynamic Regret Analysis of OFW in OCO-M

**Theorem 5** (Dynamic Regret Bound of **OFW** with Loss Variation Dependent Step Size). *Under Assumption 1 to Assumption 5, running **OFW** over unary functions $\widetilde{f}_1, \ldots, \widetilde{f}_T$ with step size $\eta = \mathcal{O}\left(\sqrt{(1+V_T)/T}\right)$ achieves against any sequence of comparators $(\mathbf{v}_1, \ldots, \mathbf{v}_T) \in \mathcal{X}^T$ the dynamic regret in eq. (9)*

$$\begin{aligned} \mathrm{Regret}_T^D &\leq \mathcal{O}\left(\sqrt{T(1+V_T+D_T)} + C_T\right) \\ &\leq \mathcal{O}\left(\sqrt{T(1+V_T+D_T+C_T)}\right). \end{aligned} \tag{48}$$

*Proof.* From the dynamic regret analysis in eq. (30), we know that running **OFW** over unary function gives

$$\sum_{t=1}^{T} \widetilde{f}_t (\mathbf{x}_t) - \sum_{t=1}^{T} \widetilde{f}_t (\mathbf{v}_t) \leq \frac{1}{\eta}(V_T + \alpha) + D\sqrt{TD_T} + \rho, \tag{49}$$

where $\alpha \triangleq 2(a + c)$ and $\rho \triangleq 4(a + c)$.

Next, we consider the switching cost of the decisions, *i.e.*, $\sum_{t=2}^{T} \|\mathbf{x}_{t-1} - \mathbf{x}_t\|_2$. By the update rule of **OFW**, we can derive an upper bound for the switching cost,

$$\sum_{t=1}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2 = \eta \sum_{t=1}^{T} \|\mathbf{x}'_{t-1} - \mathbf{x}_{t-1}\|_2 \leq \eta TD. \tag{50}$$

Combining eqs. (49) and (50), and given the definition of path length $C_T \triangleq \sum_{t=2}^{T} \|\mathbf{v}_t - \mathbf{v}_{t-1}\|_2$, we have

$$\text{Regret}_T^D \leq \eta \lambda TD + \frac{1}{\eta}(V_T + \alpha) + D\sqrt{TD_T} + \lambda C_T + \rho. \tag{51}$$

Substituting $\eta = \sqrt{(1 + V_T)/T}$ into the above equation, we directly obtain

$$\begin{aligned}
\text{Regret}_T^D &\leq \mathcal{O}\left(\sqrt{T(1 + V_T)} + \sqrt{TD_T} + C_T\right) \\
&\leq \mathcal{O}\left(\sqrt{T(1 + V_T + D_T)} + C_T\right) \\
&\leq \mathcal{O}\left(\sqrt{T(1 + V_T + D_T) + C_T^2}\right) \\
&\leq \mathcal{O}\left(\sqrt{T(1 + V_T + D_T) + TC_T}\right) \\
&= \mathcal{O}\left(\sqrt{T(1 + V_T + D_T + C_T)}\right),
\end{aligned} \tag{52}$$

where the second and third inequalities hold due to the Cauchy-Schwarz inequality, and the fourth inequality holds due to Assumption 1, *i.e.*, $0 \leq C_T = \sum_{t=2}^{T} \|\mathbf{v}_t - \mathbf{v}_{t-1}\|_2 \leq TD$. □

## C. Dynamic Regret Analysis of Meta-OFW

### C.1. Preliminaries: Online Mirror Descent

We present useful results of Online Mirror Descent (**OMD**), which enables the dynamic regret analysis for meta-algorithm, *i.e.*,**Hedge**. Consider the standard OCO setting, and the sequence of online convex functions are $\{h_t\}_{t=1,\ldots,T}$ with $h_t : \mathcal{X} \mapsto \mathbb{R}$. **OMD** starts from any $\mathbf{x}_1 \in \mathcal{X}$, and at iteration $t$, the **OMD** algorithm performs the following update

$$\mathbf{x}_{t+1} = \underset{x \in X}{argmin} \quad \eta \langle \nabla h_t (\mathbf{x}_t), \mathbf{x} \rangle + \mathcal{D}_\psi (\mathbf{x}, \mathbf{x}_t), \tag{53}$$

where $\eta > 0$ is the step size. The regularizer $\psi : \mathcal{X} \mapsto \mathbb{R}$ is a differentiable convex function defined on $\mathcal{X}$ and is assumed (without loss of generality) to be 1-strongly convex w.r.t. some norm $\|\cdot\|$ over $\mathcal{X}$. The induced Bregman divergence $\mathcal{D}_\psi$ is defined by $\mathcal{D}_\psi(\mathbf{x}, \mathbf{y}) = \psi(\mathbf{x}) - \psi(\mathbf{y}) - \langle \nabla \psi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$.

The following generic result gives an upper bound of dynamic regret with switching cost of **OMD**, which can be regarded as a generalization of OGD from gradient descent (for Euclidean norm) to mirror descent (for general primal-dual norm).

**Theorem 6.** *(Dynamic Regret Bound of **OMD** with Switching Cost (Zhao et al., 2022, Theorem 9); (Zhao et al., 2020, Theorem 2)) Provided that $\mathcal{D}_\psi(\mathbf{x}, \mathbf{z}) - \mathcal{D}_\psi(\mathbf{y}, \mathbf{z}) \leq \gamma \|\mathbf{x} - \mathbf{y}\|$ for any $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{X}$, **OMD** in eq. (53) achieves against any sequence of comparators $(\mathbf{v}_1, \ldots, \mathbf{v}_T) \in \mathcal{X}^T$ that*

$$\sum_{t=1}^{T} h_t (\mathbf{x}_t) - \sum_{t=1}^{T} h_t (\mathbf{v}_t) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\| \leq \frac{1}{\eta} \left(R^2 + \gamma C_T\right) + \eta \left(\lambda G + G^2\right) T, \tag{54}$$

*where $R^2 \triangleq \sup_{\mathbf{x},\mathbf{y} \in \mathcal{X}} \mathcal{D}_\psi(\mathbf{x}, \mathbf{y})$, $G \triangleq \sup_{\forall t} \|\nabla h_t(\cdot)\|_*$, and $\|\cdot\|_*$ is the dual norm, and $\lambda$ is a positive constant term.*

**Remark 6.** *Theorem 6 provides a way to analysis the dynamic regret and switching cost of* **OMD** *algorithm. By flexibly choosing the regularizer $\psi$ and comparator sequence $\mathbf{v}_1, \ldots, \mathbf{v}_T$, we can obtain the following two corollary (Zhao et al., 2022), which correspond to dynamic regret with switching cost of* **OGD** *(Corollary 1) and static regret with switching cost of* **Hedge** *(meta-regret) (Corollary 2), respectively.*

**Corollary 1.** *Setting the $\ell_2$ regularizer $\psi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_2^2$ and step size $\eta > 0$ for* **OMD***, suppose $\left\|\nabla \widetilde{f}_t(\mathbf{x})\right\|_2 \leq G$ and $\|\mathbf{x} - \mathbf{y}\|_2 \leq D$ hold for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ and $t \in \{1, \ldots, T\}$, then we have*

$$\sum_{t=1}^{T} \tilde{f}_t(\mathbf{x}_t) - \sum_{t=1}^{T} \widetilde{f}_t(\mathbf{v}_t) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2 \leq \frac{1}{2\eta}\left(D^2 + 2DC_T\right) + \eta\left(G^2 + \lambda G\right)T, \tag{55}$$

*which holds for any comparator sequence $\mathbf{v}_1, \ldots, \mathbf{v}_T \in \mathcal{X}$, and $C_T = \sum_{t=2}^{T} \|\mathbf{v}_{t-1} - \mathbf{v}_t\|_2$ is the path-length that measures the cumulative movements of the comparator sequence.*

Further, we present a corollary regarding the static regret with switching cost for the meta-algorithm, which is essentially a specialization of **OMD** algorithm by setting the negative-entropy regularizer.

**Corollary 2.** *Setting the negative-entropy regularizer $\psi(\boldsymbol{p}) = \sum_{i=1}^{N} p_i \log p_i$ and learning rate $\varepsilon > 0$ for* **OMD***, suppose $\|\ell_t\|_\infty \leq G$ holds for any $t \in \{1, \ldots, T\}$ and the algorithm starts from the initial weight $p_1 \in \Delta_N$, then we have*

$$\sum_{t=1}^{T} \langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle - \sum_{t=1}^{T} \ell_{t,i} + \lambda \sum_{t=2}^{T} \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1 \leq \frac{\ln(1/p_{1,i})}{\varepsilon} + \varepsilon\left(\lambda G + G^2\right)T. \tag{56}$$

Before presenting the proof of Theorem 6, we first present three useful lemmas.

**Lemma 2.** *(Chen & Teboulle, 1993, Lemma 3.2) Let $\mathcal{X}$ be a convex set in a Banach space $\mathcal{B}$. Let $f : \mathcal{X} \mapsto \mathbb{R}$ be a closed proper convex function on $\mathcal{X}$. Given a convex regularizer $\psi : \mathcal{X} \mapsto \mathbb{R}$, we denote its induced Bregman divergence by $\mathcal{D}_\psi(\cdot, \cdot)$. Then, any update of the form*

$$\mathbf{x}_k = \underset{x \in X}{\arg\min}\left\{f(\mathbf{x}) + \mathcal{D}_\psi(\mathbf{x}, \mathbf{x}_{k-1})\right\} \tag{57}$$

*satisfies the following inequality for any $\mathbf{u} \in \mathcal{X}$,*

$$f(\mathbf{x}_k) - f(\mathbf{u}) \leq \mathcal{D}_\psi(\mathbf{u}, \mathbf{x}_{k-1}) - \mathcal{D}_\psi(\mathbf{u}, \mathbf{x}_k) - \mathcal{D}_\psi(\mathbf{x}_k, \mathbf{x}_{k-1}). \tag{58}$$

**Lemma 3.** *((Duchi et al., 2010, Lemma 1);(Vishnoi, 2021)) If the regularizer $\psi : \mathcal{X} \mapsto \mathbb{R}$ is $\lambda$-strongly convex with respect to a norm $\|\cdot\|$, then we have the following lower bound for the induced Bregman divergence: $\mathcal{D}_\psi(\mathbf{x}, \mathbf{y}) \geq \frac{\lambda}{2}\|\mathbf{x} - \mathbf{y}\|^2$.*

**Lemma 4.** *(Switching Cost of* **OMD** *(Zhao et al., 2022, Lemma 10)) For* **OMD** *in eq. (53), the instantaneous switching cost is at most*

$$\|\mathbf{x}_t - \mathbf{x}_{t+1}\| \leq \eta\|\nabla h_t(\mathbf{x}_t)\|_*. \tag{59}$$

**Remark 7.** *There is an earlier result in (Zhao et al., 2020, Lemma 3) that establishes $\|\mathbf{x}_t - \mathbf{x}_{t+1}\| \leq 2\eta\|\nabla h_t(\mathbf{x}_t)\|_*$ for the switching cost of* **OMD***, instead of the inequality in eq. (59) where the multiplicative factor 2 is instead absent.*

Based on Lemma 4, we can now prove Theorem 6.

*Proof of Theorem 6.* The following proof is given in (Zhao et al., 2022, Theorem 9). The dynamic regret of **OMD** with switching cost can be bounded by following (Zhao et al., 2020, Theorem 2). The differences of (Zhao et al., 2022, Theorem 9) and (Zhao et al., 2020, Theorem 2) are that: i) (Zhao et al., 2020, Theorem 2) bound the switching cost by $\|\mathbf{x}_t - \mathbf{x}_{t+1}\| \leq 2\eta\|\nabla h_t(\mathbf{x}_t)\|_*$, instead of Lemma 4 where the multiplicative factor 2 is absent; ii) (Zhao et al., 2020, Theorem 2) derives a tighter bound on term (b) in eqs. (60) and (62) using Lemma 3.

Notice that the dynamic regret can be decomposed in the following way:

$$\sum_{t=1}^{T} h_t(\mathbf{x}_t) - \sum_{t=1}^{T} h_t(\mathbf{v}_t) \leq \sum_{t=1}^{T} \langle \nabla h_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \rangle$$

$$= \underbrace{\sum_{t=1}^{T} \langle \nabla h_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}_{t+1} \rangle}_{\text{term (a)}} + \underbrace{\sum_{t=1}^{T} \langle \nabla h_t(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{v}_t \rangle}_{\text{term (b)}}. \tag{60}$$

From Hölder's inequality and Lemma 4, we can bound term (a) as

$$\text{term (a)} \leq \sum_{t=1}^{T} \|\nabla h_t(\mathbf{x}_t)\|_* \|\mathbf{x}_t - \mathbf{x}_{t+1}\| \leq \eta \sum_{t=1}^{T} \|\nabla h_t(\mathbf{x}_t)\|_*^2. \tag{61}$$

For term (b), we have

$$
\begin{aligned}
\text{term (b)} &\leq \frac{1}{\eta} \sum_{t=1}^{T} \left( \mathcal{D}_\psi(\mathbf{v}_t, \mathbf{x}_t) - \mathcal{D}_\psi(\mathbf{v}_t, \mathbf{x}_{t+1}) - \mathcal{D}_\psi(\mathbf{x}_{t+1}, \mathbf{x}_t) \right) \\
&\leq \frac{1}{\eta} \sum_{t=2}^{T} \left( \mathcal{D}_\psi(\mathbf{v}_t, \mathbf{x}_t) - \mathcal{D}_\psi(\mathbf{v}_{t-1}, \mathbf{x}_t) \right) + \frac{1}{\eta} \mathcal{D}_\psi(\mathbf{v}_1, \mathbf{x}_1) \\
&\leq \frac{\gamma}{\eta} \sum_{t=2}^{T} \|\mathbf{v}_t - \mathbf{v}_{t-1}\| + \frac{1}{\eta} R^2,
\end{aligned}
\tag{62}
$$

where the first inequality holds by Lemma 2, the second inequality holds by the non-negativity of the Bregman divergence, and the last inequality holds due to $\mathcal{D}_\psi(\mathbf{x}, \mathbf{z}) - \mathcal{D}_\psi(\mathbf{y}, \mathbf{z}) \leq \gamma \|\mathbf{x} - \mathbf{y}\|$ for any $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{X}$.

By Lemma 3, the switching cost is bounded as

$$\sum_{t=2}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\| \leq \eta \sum_{t=2}^{T} \|\nabla h_{t-1}(\mathbf{x}_{t-1})\|_*. \tag{63}$$

Combining eqs. (61) to (63), we obtain

$$
\begin{aligned}
\sum_{t=1}^{T} h_t(\mathbf{x}_t) - \sum_{t=1}^{T} h_t(\mathbf{v}_t) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\| &\leq \frac{1}{\eta} \left( R^2 + \gamma C_T \right) + \eta \sum_{t=1}^{T} \left( \lambda \|\nabla h_t(\mathbf{x}_t)\|_* + \|\nabla h_{t-1}(\mathbf{x}_{t-1})\|_*^2 \right) \\
&\leq \frac{1}{\eta} \left( R^2 + \gamma C_T \right) + \eta \left( \lambda G + G^2 \right) T. \quad \square
\end{aligned}
\tag{64}
$$

## C.2. Proof of Theorem 2

*Proof.* According to (Anava et al., 2015, Proof of Theorem 3.1), the coordinate-Lipschitz continuity of $f_t$ (Assumption 4) implies that

$$
\begin{aligned}
\left| f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t) - \widetilde{f}_t(\mathbf{x}_t) \right| &\leq L \sum_{i=1}^{m} \|\mathbf{x}_t - \mathbf{x}_{t-i}\|_2 \\
&\leq L \sum_{i=1}^{m} \sum_{l=1}^{i} \|\mathbf{x}_{t-l+1} - \mathbf{x}_{t-l}\|_2 \\
&\leq mL \sum_{i=1}^{m} \|\mathbf{x}_{t-i+1} - \mathbf{x}_{t-i}\|_2.
\end{aligned}
\tag{65}
$$

Taking the summation from $t = 1$ to $T$ gives

$$\left| \sum_{t=1}^{T} f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t) - \sum_{t=1}^{T} \widetilde{f}_t(\mathbf{x}_t) \right| \leq mL \sum_{t=1}^{T} \sum_{i=1}^{m} \|\mathbf{x}_{t-i+1} - \mathbf{x}_{t-i}\|_2 \leq m^2 L \sum_{t=1}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2, \tag{66}$$

and the dynamic regret can be thus upper bounded by

$$
\begin{aligned}
\text{Regret}_T^D &= \sum_{t=1}^{T} f_t(\mathbf{x}_{t-m}, \ldots, \mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{v}_{t-m}, \ldots, \mathbf{v}_t) \\
&\leq \underbrace{\sum_{t=1}^{T} \widetilde{f}_t(\mathbf{x}_t) - \sum_{t=1}^{T} \widetilde{f}_t(\mathbf{v}_t)}_{\text{dynamic regret over unary loss}} + \underbrace{\lambda \sum_{t=1}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2}_{\text{switching cost of decisions}} + \underbrace{\lambda \sum_{t=1}^{T} \|\mathbf{v}_t - \mathbf{v}_{t-1}\|_2}_{\text{switching cost of comparators}},
\end{aligned}
\tag{67}
$$

where $\mathbf{x}_\tau$ and $\mathbf{v}_\tau$ can be set as $\mathbf{0}$ for all $\tau \leq 0$, and $\lambda \triangleq m^2 L$.

By Lemma 5, we aim to bound the meta-regret and base-regret terms.

**Meta-regret bound.** Denote by $\mathbf{e}_i$ the $i$-th standard basis of $\mathbb{R}^N$-space. Since the meta-algorithm performs Hedge over the switching-cost-regularized loss $\boldsymbol{\ell}_t \in \mathbb{R}^N$, Corollary 2 implies that for any $i \in \{1, \ldots, N\}$,

$$
\begin{aligned}
\sum_{t=1}^{T} \langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle - \sum_{t=1}^{T} \ell_{t,i} + \lambda D \sum_{t=2}^{T} \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1 &\leq \varepsilon \left( \lambda D G_{\text{meta}} + G_{\text{meta}}^2 \right) T + \frac{\mathcal{D}_\psi (\mathbf{e}_i, \boldsymbol{p}_1)}{\varepsilon} \\
&= \varepsilon (2\lambda + G)(\lambda + G)D^2 T + \frac{\ln (1/p_{1,i})}{\varepsilon} \\
&\leq \varepsilon (2\lambda + G)(\lambda + G)D^2 T + \frac{2 \ln(i+1)}{\varepsilon}.
\end{aligned}
\tag{68}
$$

where the first inequality holds due to $G_{\text{meta}} = \max_{t \in \{1,\ldots,T\}} \|\boldsymbol{\ell}_t\|_\infty \leq (\lambda + G)D$, and the last inequality holds by plugging in the initialization of weights *i.e.*, $\boldsymbol{p}_1 \in \Delta_N$ with $p_{1,i} = \frac{1}{i(i+1)} \cdot \frac{N+1}{N}$ for any $i \in \{1, \ldots, N\}$. By choosing the optimal learning rate $\varepsilon = \varepsilon^* = \sqrt{\frac{2}{(2\lambda+G)(\lambda+G)D^2 T}}$, we can obtain the following upper bound for the meta-regret,

$$
\sum_{t=1}^{T} \langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle - \sum_{t=1}^{T} \ell_{t,i} + \lambda D \sum_{t=2}^{T} \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1 \leq D\sqrt{2(2\lambda + G)(\lambda + G)T}(1 + \ln(i+1)).
\tag{69}
$$

**Base-regret bound.** As specified by **Meta-OFW**, there are multiple base-learners, each performing **OFW** over the linearized loss with a particular step size $\eta_i \in \mathcal{H}$ for base-learner $\mathcal{B}_i$. As a result, eqs. (49) and (50) and the definition of $\bar{D}_T$ imply that the base-regret satisfies

$$
\sum_{t=1}^{T} g_t (\mathbf{x}_{t,i}) - \sum_{t=1}^{T} g_t (\mathbf{v}_t) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2 \leq \eta_i \lambda T D + \frac{1}{\eta_i}(V_T + \alpha) + D\sqrt{T \bar{D}_T} + \rho,
\tag{70}
$$

which holds for any comparator sequence $\mathbf{v}_1, \ldots, \mathbf{v}_T \in \mathcal{X}$ as well as any base-learner $i \in \{1, \ldots, N\}$.

**Dynamic regret bound.** Due to the boundedness of the loss variation $V_T$, we know that the optimal step size $\eta_*$ provably lies in the range of $[\eta_1, \eta_N]$. In particular, given eq. (70), the optimal step size $\eta_*$ is

$$
\sqrt{\frac{\alpha}{\lambda T D}} \leq \eta_* = \sqrt{\frac{V_T + \alpha}{\lambda T D}} \leq \sqrt{\frac{Tc + \alpha}{\lambda T D}}.
\tag{71}
$$

Furthermore, by the construction of the step size pool in eq. (12), there exists an index $i^* \in \{1, \ldots, N\}$, such that $\eta_{i^*} \leq \eta_* \leq \eta_{i^*+1} = 2\eta_{i^*}$, with

$$
i^* \leq \left\lceil \frac{1}{2} \log_2 \left( 1 + \frac{V_T}{\alpha} \right) \right\rceil + 1.
\tag{72}
$$

Notice that the meta-base decomposition in Lemma 5 holds for any index of base-learners $i$. Thus, in particular, we can

choose the index $i^*$ and achieve the following result by using the meta-regret and base-regret bounds in eqs. (69) and (70),

$$
\begin{aligned}
&\sum_{t=1}^{T} \widetilde{f}_t\left(\mathbf{x}_t\right) - \sum_{t=1}^{T} \tilde{f}_t\left(\mathbf{v}_t\right) + \lambda \sum_{t=2}^{T} \left\|\mathbf{x}_t - \mathbf{x}_{t-1}\right\|_2 \\
&\leq \underbrace{\sum_{t=1}^{T}\left(\langle\boldsymbol{p}_t, \boldsymbol{\ell}_t\rangle - \ell_{t, i^*}\right) + \lambda D \sum_{t=2}^{T}\left\|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\right\|_1}_{\text{meta-regret}} + \underbrace{\sum_{t=1}^{T}\left(g_t\left(\mathbf{x}_{t, i^*}\right) - g_t\left(\mathbf{v}_t\right)\right) + \lambda \sum_{t=2}^{T}\left\|\mathbf{x}_{t, i^*} - \mathbf{x}_{t-1, i^*}\right\|_2}_{\text{base-regret}} \\
&\leq D \sqrt{2(2\lambda+G)(\lambda+G) T}\left(1+\ln\left(i^*+1\right)\right) + \left(\eta_{i^*} \lambda T D + \frac{1}{\eta_{i^*}}\left(V_T+\alpha\right) + D \sqrt{T \bar{D}_T} + \rho\right) \\
&\leq D \sqrt{2(2\lambda+G)(\lambda+G) T}\left(1+\ln\left(i^*+1\right)\right) + \eta_* \lambda T D + \frac{2}{\eta_*}\left(V_T+\alpha\right) + D \sqrt{T \bar{D}_T} + \rho \\
&\leq \underbrace{2 D(\lambda+G) \sqrt{T}\left(1+\ln\left(\left\lceil\frac{1}{2} \log_2\left(1+\frac{V_T}{\alpha}\right)\right\rceil+2\right)\right)}_{\leq \mathcal{O}\left(\sqrt{T}\left(1+\log \log V_T\right)\right)} + \underbrace{3 \sqrt{\lambda T D\left(V_T+\alpha\right)} + D \sqrt{T \bar{D}_T}}_{\leq \mathcal{O}\left(\sqrt{T\left(1+V_T+\bar{D}_T\right)}\right)} + \rho \\
&\leq \mathcal{O}\left(\sqrt{T\left(1+V_T+\bar{D}_T\right)}\right).
\end{aligned}
\tag{73}
$$

Combining eq. (9) and eq. (73) gives

$$
\begin{aligned}
\operatorname{Regret}_T^D &\leq \sum_{t=1}^{T} \widetilde{f}_t\left(\mathbf{x}_t\right) - \sum_{t=1}^{T} \widetilde{f}_t\left(\mathbf{v}_t\right) + \lambda \sum_{t=2}^{T}\left\|\mathbf{x}_t - \mathbf{x}_{t-1}\right\|_2 + \lambda \sum_{t=2}^{T}\left\|\mathbf{v}_t - \mathbf{v}_{t-1}\right\|_2 \\
&\leq \mathcal{O}\left(\sqrt{T\left(1+V_T+\bar{D}_T\right)}\right) + \mathcal{O}\left(C_T\right) \\
&\leq \mathcal{O}\left(\sqrt{T\left(1+V_T+\bar{D}_T\right) + C_T^2}\right) \\
&\leq \mathcal{O}\left(\sqrt{T\left(1+V_T+\bar{D}_T\right) + T C_T}\right) \\
&= \mathcal{O}\left(\sqrt{T\left(1+V_T+\bar{D}_T+C_T\right)}\right),
\end{aligned}
\tag{74}
$$

where the third inequality holds due to Cauchy-Schwartz inequality, and the fourth inequality holds due to Assumption 1, i.e., $0 \leq C_T = \sum_{t=2}^{T}\left\|\mathbf{v}_t - \mathbf{v}_{t-1}\right\|_2 \leq T D$. $\qquad\square$

### C.3. Decomposition of Unary Cost and Switching Cost (Zhao et al., 2022)

**Lemma 5** (Decomposition of Unary Cost and Switching Cost (Zhao et al., 2022)). *The unary and switching costs in eq. (9) can be decomposed as follows:*

$$
\begin{aligned}
\sum_{t=1}^{T} \tilde{f}_t\left(\mathbf{x}_t\right) - \sum_{t=1}^{T} \tilde{f}_t\left(\mathbf{v}_t\right) + \lambda \sum_{t=2}^{T}\left\|\mathbf{x}_t - \mathbf{x}_{t-1}\right\|_2 &\leq \underbrace{\sum_{t=1}^{T}\left(\langle\boldsymbol{p}_t, \boldsymbol{\ell}_t\rangle - \ell_{t, i}\right) + \lambda D \sum_{t=2}^{T}\left\|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\right\|_1}_{\text{meta-regret}} \\
&+ \underbrace{\sum_{t=1}^{T}\left(g_t\left(\mathbf{x}_{t, i}\right) - g_t\left(\mathbf{v}_t\right)\right) + \lambda \sum_{t=2}^{T}\left\|\mathbf{x}_{t, i} - \mathbf{x}_{t-1, i}\right\|_2}_{\text{base-regret}},
\end{aligned}
\tag{75}
$$

*which holds for any $i \in \{1, \ldots, N\}$.*

*Proof.* By the meta-base structure, the final decision of each round is $\mathbf{x}_t = \sum_{i=1}^{N} p_{t, i} \mathbf{x}_{t, i}$. Therefore, we can expand the

switching cost of the final prediction sequence as

$$
\begin{aligned}
\|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2 &= \left\| \sum_{i=1}^{N} p_{t,i} \mathbf{x}_{t,i} - \sum_{i=1}^{N} p_{t-1,i} \mathbf{x}_{t-1,i} \right\|_2 \\
&\leq \left\| \sum_{i=1}^{N} p_{t,i} \mathbf{x}_{t,i} - \sum_{i=1}^{N} p_{t,i} \mathbf{x}_{t-1,i} \right\|_2 + \left\| \sum_{i=1}^{N} p_{t,i} \mathbf{x}_{t-1,i} - \sum_{i=1}^{N} p_{t-1,i} \mathbf{x}_{t-1,i} \right\|_2 \\
&\leq \sum_{i=1}^{N} p_{t,i} \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2 + D \sum_{i=1}^{N} |p_{t,i} - p_{t-1,i}| \\
&= \sum_{i=1}^{N} p_{t,i} \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2 + D \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1 ,
\end{aligned}
\tag{76}
$$

where the first inequality holds due to the triangle inequality and the second inequality is true owing to the boundedness of the feasible domain (Assumption 1).

By eq. (76) and convexity of $\widetilde{f}_t(\cdot)$, we have

$$
\begin{aligned}
&\sum_{t=1}^{T} \widetilde{f}_t(\mathbf{x}_t) - \sum_{t=1}^{T} \widetilde{f}_t(\mathbf{v}_t) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2 \\
&\leq \sum_{t=1}^{T} \left\langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \right\rangle + \lambda D \sum_{t=2}^{T} \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1 + \lambda \sum_{t=2}^{T} \sum_{i=1}^{N} p_{t,i} \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2 .
\end{aligned}
\tag{77}
$$

Next, we manipulate eq. (77) using the definition of linearized loss $g_t(\mathbf{x}) = \left\langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{x} \right\rangle$,

$$
\begin{aligned}
&\sum_{t=1}^{T} \widetilde{f}_t(\mathbf{x}_t) - \sum_{t=1}^{T} \widetilde{f}_t(\mathbf{v}_t) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2 \\
&\leq \sum_{t=1}^{T} \sum_{i=1}^{N} p_{t,i} \left( \left\langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{x}_{t,i} \right\rangle + \lambda \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2 \right) - \sum_{t=1}^{T} \left( \left\langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{x}_{t,i} \right\rangle + \lambda \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2 \right) \\
&\quad + \lambda D \sum_{t=2}^{T} \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1 + \sum_{t=1}^{T} \left( \left\langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{x}_{t,i} \right\rangle - \left\langle \nabla \widetilde{f}_t(\mathbf{x}_t), \mathbf{v}_t \right\rangle \right) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2 \\
&= \underbrace{\sum_{t=1}^{T} (\langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle - \ell_{t,i}) + \lambda D \sum_{t=2}^{T} \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1}_{\text{meta-regret}} + \underbrace{\sum_{t=1}^{T} (g_t(\mathbf{x}_{t,i}) - g_t(\mathbf{v}_t)) + \lambda \sum_{t=2}^{T} \|\mathbf{x}_{t,i} - \mathbf{x}_{t-1,i}\|_2}_{\text{base-regret}} . \quad \square
\end{aligned}
\tag{78}
$$

### C.4. Proof of Corollary 2 (Zhao et al., 2022)

*Proof.* From the proof of Theorem 6, we can obtain that

$$
\sum_{t=1}^{T} \langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle - \sum_{t=1}^{T} \ell_{t,i} + \lambda \sum_{t=2}^{T} \|\boldsymbol{p}_t - \boldsymbol{p}_{t-1}\|_1 \leq \frac{\mathcal{D}_\psi(\mathbf{e}_i, \boldsymbol{p}_1)}{\varepsilon} + \varepsilon \left( \lambda G + G^2 \right) T,
\tag{79}
$$

where $\mathbf{e}_i$ the $i$-th standard basis of $\mathbb{R}^N$-space. When choosing the negative-entropy regularizer, the induced Bregman divergence becomes Kullback-Leibler divergence, i.e., $\mathcal{D}_\psi(\boldsymbol{q}, \boldsymbol{p}) = \mathrm{KL}(\boldsymbol{q}, \boldsymbol{p}) = \sum_{i=1}^{N} q_i \ln(q_i/p_i)$. Therefore, $\mathcal{D}_\psi(\boldsymbol{e}_i, \boldsymbol{p}_1) = \ln(1/p_{1,i})$, which implies the desired result. $\square$

## D. Dynamic Regret Analysis of Non-Stochastic Control

### D.1. Definition of Strongly Stable Controller

**Definition 5.** *A linear policy $K$ is $(\kappa, \gamma)$-strongly stable if there exist matrices $L$ and $Q$ satisfying $A - BK = QLQ^{-1}$, such that following two conditions are satisfied:*

1. *The spectral norm of L is strictly smaller than one,* i.e., $\|L\|_{\mathrm{op}} \leq 1 - \gamma$;

2. *The controller and the transforming matrices are bounded,* $\|K\|_{\mathrm{op}}$, $\|Q\|_{\mathrm{op}}$, *and* $\|Q\|_{\mathrm{op}}^{-1} \leq \kappa$.

### D.2. Proof of Theorem 3

*Proof.* We decompose the dynamic regret as follows,

$$
\sum_{t=0}^{T} c_t\left(x_t, u_t\right) - \sum_{t=0}^{T} c_t\left(x_t^*, u_t^*\right)
$$

$$
= \sum_{t=0}^{T} c_t\left(x_t\left(M_{0:t-1}\right), u_t\left(M_{0:t}\right)\right) - \sum_{t=0}^{T} c_t\left(x_t^*\left(0\right), u_t^*\left(0\right)\right)
$$

$$
= \underbrace{\sum_{t=0}^{T} c_t\left(x_t\left(M_{0:t-1}\right), u_t\left(M_{0:t}\right)\right) - \sum_{t=0}^{T} f_t\left(M_{t-1-H:t}\right)}_{\text{term (a)}} + \underbrace{\sum_{t=0}^{T} f_t\left(M_{t-1-H:t}\right) - \sum_{t=0}^{T} f_t\left(M_{t-1-H:t}^*\right)}_{\text{term (b)}}
$$

$$
+ \underbrace{\sum_{t=0}^{T} f_t\left(M_{t-1-H:t}^*\right) - \sum_{t=0}^{T} c_t\left(x_t\left(M_{0:t-1}^*\right), u_t\left(M_{0:t}^*\right)\right)}_{\text{term (c)}} + \underbrace{\sum_{t=0}^{T} c_t\left(x_t\left(M_{0:t-1}^*\right), u_t\left(M_{0:t}^*\right)\right) - \sum_{t=0}^{T} c_t\left(x_t^*\left(0\right), u_t^*\left(0\right)\right)}_{\text{term (d)}},
$$

$$\tag{80}$$

where term (a) and term (c) are the approximation errors induced by truncation of loss function, term (b) is the dynamic regret over the truncated loss functions $\{f_t\}_{t=0,\ldots,T}$, and term (d) is the approximation error of DAC controller.

By Theorem 8, term (a) and term (c) are bounded by

$$
\text{term (a)} + \text{term (c)} \leq 4TG_c D^2 \kappa^3 (1-\gamma)^{H+1}. \tag{81}
$$

Then by Proposition 2, we can bound term (d) as

$$
\text{term (d)} \leq \frac{4TG_c DWH\kappa_B^2\kappa^6(1-\gamma)^{H-1}}{\gamma}. \tag{82}
$$

Next, we focus on the term (b). Similar to eq. (9), we decompose term (b) as follows,

$$
\text{term (b)} = \sum_{t=0}^{T} f_t\left(M_{t-1-H:t}\right) - \sum_{t=0}^{T} f_t\left(M_{t-1-H:t}^*\right)
$$

$$
\leq \sum_{t=0}^{T} \widetilde{f}_t\left(M_t\right) - \sum_{t=0}^{T} \widetilde{f}_t\left(M_t^*\right) + \zeta\sum_{t=1}^{T} \|M_{t-1} - M_t\|_{\mathrm{F}} + \zeta\sum_{t=1}^{T} \|M_{t-1}^* - M_t^*\|_{\mathrm{F}}
$$

$$
\leq \sum_{t=0}^{T} \left\langle \nabla_M \widetilde{f}_t\left(M_t\right), M_t - M_t^* \right\rangle + \zeta\sum_{t=1}^{T} \|M_{t-1} - M_t\|_{\mathrm{F}} + \zeta\sum_{t=1}^{T} \|M_{t-1}^* - M_t^*\|_{\mathrm{F}} \tag{83}
$$

$$
= \underbrace{\sum_{t=0}^{T} g_t\left(M_t\right) - \sum_{t=0}^{T} g_t\left(M_t^*\right) + \zeta\sum_{t=1}^{T} \|M_{t-1} - M_t\|_{\mathrm{F}}}_{\text{Dynamic Regret with Switching Cost over } \{g_t\}_{t=0,\ldots,T}} + \underbrace{\zeta\sum_{t=1}^{T} \|M_{t-1}^* - M_t^*\|_{\mathrm{F}}}_{\text{Path Length of Comparators}},
$$

where $\zeta = (H+2)^2 L_f$ and $g_t(M) = \left\langle \nabla_M \widetilde{f}_t\left(M_t\right), M \right\rangle$ is the surrogate linearized loss.

We now aim to analyze the dynamic regret with switching cost in eq. (83), which can be decomposed into meta-regret and

base-regret similar to Lemma 5:

$$
\begin{aligned}
&\sum_{t=0}^{T} g_t\left(M_t\right) - \sum_{t=0}^{T} g_t\left(M_t^*\right) + \zeta \sum_{t=1}^{T} \left\|M_{t-1} - M_t\right\|_{\mathrm{F}} \\
&= \underbrace{\sum_{t=0}^{T} \left\langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \right\rangle - \sum_{t=0}^{T} \ell_{t,i} + \zeta D_f \sum_{t=1}^{T} \left\|\boldsymbol{p}_{t-1} - \boldsymbol{p}_t\right\|_1}_{\text{meta-regret}} + \underbrace{\left( \sum_{t=0}^{T} g_t\left(M_{t,i}\right) - \sum_{t=0}^{T} g_t\left(M_t^*\right) \right) + \zeta \sum_{t=1}^{T} \left\|M_{t-1,i} - M_{t,i}\right\|_{\mathrm{F}}}_{\text{base-regret}},
\end{aligned}
\tag{84}
$$

where $\ell_t \in \Delta_N$ is the surrogate loss vector of the meta-algorithm with $\ell_{t,i} = \zeta \left\|M_{t-1,i} - M_{t,i}\right\|_{\mathrm{F}} + g_t\left(M_{t,i}\right)$, for $i \in \{1, \ldots, N\}$. Note that the regret decomposition holds for any base-learner $\mathcal{B}_i$.

**Meta-regret bound.** By Corollary 2, we obtain

$$
\begin{aligned}
\text{meta-regret} &\leq \varepsilon \left(2\zeta + G_f\right)\left(\zeta_f + G_f\right) D_f^2(T+1) + \frac{\ln\left(1/p_{1,i}\right)}{\varepsilon} \\
&= D_f \sqrt{2\left(2\zeta + G_f\right)\left(\zeta + G_f\right)(T+1)}(1 + \ln(1+i)),
\end{aligned}
\tag{85}
$$

where the equality holds by choosing the optimal learning rate $\varepsilon = \varepsilon^* = \sqrt{\frac{2}{D_f^2(2\zeta + G_f)(\zeta + G_f)(T+1)}}$.

**Base-regret bound.** By Theorem 7 and the definition of $\bar{D}_T$, we have

$$
\begin{aligned}
\text{base-regret} &\leq \eta_i \zeta T D_f + \frac{1}{\eta_i}(V_T + \sigma) + D_f \sqrt{(T+1)D_T} + \theta \\
&\leq \eta_i \zeta T D_f + \frac{1}{\eta_i}(V_T + \sigma) + D_f \sqrt{(T+1)\bar{D}_T} + \theta,
\end{aligned}
\tag{86}
$$

where $V_T \triangleq \sum_{t=0}^{T} \sup_{M \in \mathcal{M}} \left|f_t(M) - f_{t-1}(M)\right|$, $D_T \triangleq \sum_{t=0}^{T} \left\|\nabla_M f_t\left(M_t\right) - \nabla_M f_{t-1}\left(M_{t-1}\right)\right\|_{\mathrm{F}}^2$, $\sigma \triangleq 4\beta D^2$, and $\theta \triangleq 8\beta D^2$.

**Overall regret bound.** Due to the boundedness of the loss variation $V_T$, the optimal step size $\eta_*$ provably lies in the range of $[\eta_1, \eta_N]$. In particular, the optimal step size $\eta_*$ is

$$
\sqrt{\frac{\sigma}{\zeta T D_f}} \leq \eta_* = \sqrt{\frac{V_T + \sigma}{\zeta T D_f}} \leq \sqrt{\frac{2\beta D^2 T + \phi}{\zeta T D_f}},
\tag{87}
$$

where $\phi = \sigma + 2\beta D^2$.

Furthermore, by the construction of the step size pool in eq. (12), there exists an index $i^* \in \{1, \ldots, N\}$, such that $\eta_{i^*} \leq \eta_* \leq \eta_{i^*+1} = 2\eta_{i^*}$, with

$$
i^* \leq \left\lceil \frac{1}{2} \log_2 \left(1 + \frac{V_T}{\sigma}\right) \right\rceil + 1.
\tag{88}
$$

Since the meta-base decomposition in eq. (84) holds for any index $i$, we can choose the index $i^*$ and achieve the following result by using the meta-regret and base-regret bounds in eqs. (85) and (86),

$$
\begin{aligned}
&\sum_{t=0}^{T} g_t\left(M_t\right) - \sum_{t=0}^{T} g_t\left(M_t^*\right) + \zeta \sum_{t=1}^{T} \left\|M_{t-1} - M_t\right\|_{\mathrm{F}} \\
&\leq D_f \sqrt{2\left(2\zeta + G_f\right)\left(\zeta + G_f\right)(T+1)}(1 + \ln(i^* + 1)) + \left( \eta_{i^*} T D_f + \frac{1}{\eta_{i^*}}(V_T + \sigma) + D_f \sqrt{(T+1)\bar{D}_T} + \theta \right) \\
&\leq D_f \sqrt{2\left(2\zeta + G_f\right)\left(\zeta + G_f\right)(T+1)}(1 + \ln(i^* + 1)) + \eta_* T D_f + \frac{2}{\eta_*}(V_T + \sigma) + D_f \sqrt{(T+1)\bar{D}_T} + \theta \\
&\leq D_f \sqrt{2\left(2\zeta + G_f\right)\left(\zeta + G_f\right)(T+1)}(1 + \ln(i^* + 1)) \left(1 + \ln\left(\left\lceil \frac{1}{2} \log_2 \left(1 + \frac{V_T}{\sigma}\right) \right\rceil + 2\right)\right) \\
&\quad + 3\sqrt{T\left(V_T + k_3\right)} + D_f \sqrt{(T+1)\bar{D}_T} + \theta.
\end{aligned}
\tag{89}
$$

---

**Algorithm 4 OFW** for Non-Stochastic Control.

---

**Input:** Time horizon $T$; step size $\eta$.
**Output:** Prediction $M_t$ at each time step $t = 1, \ldots, T$.

1: Initialize $M_0 \in \mathcal{M}$;
2: **for** each time step $t = 1, \ldots, T$ **do**
3:     Obtain gradient $\nabla_M f_t(M_t)$;
4:     Compute $M_t' = \arg\min_{M \in \mathcal{M}} \langle \nabla_M f_t(M_t), M \rangle_{\mathrm{F}}$;
5:     Update $M_{t+1} = (1 - \eta)M_t + \eta M_t'$;
6: **end for**

---

Combining eqs. (81) to (83) and (89), and $C_T \triangleq \sum_{t=2}^{T} \left\| M_{t-1}^* - M_t^* \right\|_{\mathrm{F}}$ gives

$$
\begin{aligned}
&\sum_{t=0}^{T} c_t\left(x_t, u_t\right) - \sum_{t=0}^{T} c_t\left(x_t^*, u_t^*\right) \\
&\leq 4TG_c D^2 \kappa^3 (1-\gamma)^{H+1} + \frac{4TG_c DWH\kappa_B^2 \kappa^6 (1-\gamma)^{H-1}}{\gamma} \\
&\quad + D_f \sqrt{2\left(2\zeta + G_f\right)\left(\zeta + G_f\right)(T+1)}(1 + \ln(i^* + 1)) \left(1 + \ln\left(\left\lceil \frac{1}{2}\log_2\left(1 + \frac{V_T}{\sigma}\right)\right\rceil + 2\right)\right) \\
&\quad + 3\sqrt{T\left(V_T + k_3\right)} + D_f \sqrt{(T+1)\bar{D}_T} + \theta + \zeta C_T.
\end{aligned}
\tag{90}
$$

Finally, by setting $H = \mathcal{O}(\log T)$, the final dynamic policy regret is bounded by $\widetilde{\mathcal{O}}\left(\sqrt{T\left(1 + V_T + \bar{D}_T + C_T\right)}\right)$. $\qquad\square$

### D.3. Dynamic Regret Analysis of OFW over $\mathcal{M}$-space

In Theorem 1, we have analyzed the dynamic regret of **Meta-OFW** over the Euclidean space. To utilize **Meta-OFW** for non-stochastic control, we need to generalized the result to $\mathcal{M}$-space, *i.e.*, generalize the previous results from Euclidean norm to Frobenius norm over $\mathcal{M}$-space. To this end, we first present the dynamic regret analysis of **OFW** over $\mathcal{M}$-space.

**Theorem 7.** *Suppose the function* $\tilde{f} : \mathcal{M} \mapsto \mathbb{R}$ *is convex, with bounded the gradient norm* $G_f$, *i.e.,* $\left\| \nabla_M \widetilde{f}_t(M) \right\|_{\mathrm{F}} \leq G_f$ *for any* $M \in \mathcal{M}$ *and* $t \in \{1, \ldots, T\}$, *and bounded Euclidean diameter of* $\mathcal{M}$-space $D_f$, *i.e.,* $\sup_{M, M' \in \mathcal{M}} \|M - M'\|_{\mathrm{F}} \leq D_f$. *Then for any comparator sequence* $M_1, \ldots, M_T \in \mathcal{M}$, *Algorithm 4 satisfies that*

$$
\sum_{t=1}^{T} \widetilde{f}_t\left(M_t\right) - \sum_{t=1}^{T} \widetilde{f}_t\left(M_t^*\right) + \zeta \sum_{t=2}^{T} \left\| M_{t-1} - M_t \right\|_{\mathrm{F}} \leq \eta T \zeta D_f + \frac{1}{\eta}(V_T + \sigma) + D_f \sqrt{T D_T} + \theta,
\tag{91}
$$

*where* $V_T \triangleq \sum_{t=1}^{T} \sup_{M \in \mathcal{M}} |f_t(M) - f_{t-1}(M)|$, *and* $D_T \triangleq \sum_{t=1}^{T} \left\| \nabla_M f_t\left(M_t\right) - \nabla_M f_{t-1}\left(M_{t-1}\right) \right\|_{\mathrm{F}}^2$, $\sigma \triangleq 4\beta D^2$, *and* $\theta \triangleq 8\beta D^2$.

*Proof.* Consider the term $\sum_{t=1}^{T} \widetilde{f}_t\left(M_t\right) - \sum_{t=1}^{T} \widetilde{f}_t\left(M_t^*\right)$. The bound can be developed similar to Theorem 1, replacing vector inner product $\langle \cdot, \cdot \rangle$ by matrix inner product $\langle \cdot, \cdot \rangle_{\mathrm{F}}$ and replacing vector norm $\|\cdot, \cdot\|_2$ by Frobenius norm $\|\cdot, \cdot\|_{\mathrm{F}}$. Then from eq. (30), we directly obtain

$$
\sum_{t=1}^{T} \widetilde{f}_t\left(M_t\right) - \sum_{t=1}^{T} \widetilde{f}_t\left(M_t^*\right) \leq \frac{1}{\eta}(V_T + \sigma) + D_f \sqrt{T D_T} + \theta.
\tag{92}
$$

On the other hand, the switching cost can be bounded by

$$
\sum_{t=2}^{T} \left\| M_t - M_{t-1} \right\|_{\mathrm{F}} = \eta \sum_{t=2}^{T} \left\| M_{t-1}' - M_{t-1} \right\|_{\mathrm{F}} \leq \eta(T-1)D_f \leq \eta T D_f.
\tag{93}
$$

We complete the proof by combining the above equations. $\qquad\square$

## D.4. State Transition under DAC Controller

**Proposition 1** (State Transition under DAC Controller). *Suppose the initial state is $x_0 = 0$, and one chooses the DAC controller $\pi(M_t, K_t)$ at iteration $t$, the reaching state is*

$$x_{t+1} = \widetilde{A}_{K_{t:t-h}} x_{t-h} + \sum_{i=0}^{H+h} \Psi_{t,i}^{K_t,h}(M_{t-h:t}) w_{t-i}, \tag{94}$$

*where $\Psi_{t,i}^{K,h}(M_{t-h:t})$ is the transfer matrix defined as*

$$\Psi_{t,i}^{K_t,h}(M_{t-h:t}) = \widetilde{A}_{K_{t:t-i+1}} \mathbf{1}_{i\leq h} + \sum_{j=0}^{h} \widetilde{A}_{K_{t:t-j+1}} B_{t-j} M_{t-j}^{[i-j-1]} \mathbf{1}_{1\leq i-j\leq H}, \tag{95}$$

*where $\widetilde{A}_{K_{t:t-i}} \triangleq \prod_{\tau=t}^{t-i}(A_\tau - B_\tau K_\tau)$, and we define $\widetilde{A}_{K_{t:t-i}} \triangleq \mathbf{I}$ if $i < 0$. The evolving equation holds for any $h \in \{0, \ldots, t\}$.*

*Proof.* The proof follows the same step as in (Agarwal et al., 2019, Lemma 4.3). We aim to show

$$x_{t+1} = \widetilde{A}_{K_{t:t-h}} x_{t-h} + \sum_{i=0}^{H+h} \Psi_{t,i}^{K_t,h}(M_{t-h:t}) w_{t-i}, \tag{96}$$

where $\Psi_{t,i}^{K,h}(M_{t-h:t})$ is the transfer matrix defined as

$$\Psi_{t,i}^{K_t,h}(M_{t-h:t}) = \widetilde{A}_{K_{t:t-i+1}} \mathbf{1}_{i\leq h} + \sum_{j=0}^{h} \widetilde{A}_{K_{t:t-j+1}} B_{t-j} M_{t-j}^{[i-j-1]} \mathbf{1}_{1\leq i-j\leq H}. \tag{97}$$

For any time $t$, we will prove the claim by induction. For $h = 0$, we have

$$
\begin{aligned}
x_{t+1} &= \widetilde{A}_{K_t} x_t + \sum_{i=1}^{H} B_t M_t^{[i-1]} w_{t-i} + w_t \\
&= \widetilde{A}_{K_t} x_t + \sum_{i=0}^{H} \Psi_{t,i}^{K_t,0}(M_t) w_{t-i}.
\end{aligned}
\tag{98}
$$

Suppose that eq. (96) holds for some $h \geq 0$, then for $h + 1$ we have

$$
\begin{aligned}
x_{t+1} &= \widetilde{A}_{K_{t:t-h}} x_{t-h} + \sum_{i=0}^{H+h} \Psi_{t,i}^{K_t,h}(M_{t-h:t}) w_{t-i} \\
&= \widetilde{A}_{K_{t:t-h}} \left( \widetilde{A}_{K_{t-h-1}} x_{t-h-1} + \sum_{i=1}^{H} B_{t-h-1} M_{t-h-1}^{[i-1]} w_{t-h-1-i} + w_{t-h-1} \right) + \sum_{i=0}^{H+h} \Psi_{t,i}^{K_t,h}(M_{t-h:t}) w_{t-i} \\
&= \widetilde{A}_{K_{t:t-h-1}} x_{t-h-1} + \sum_{i=0}^{H+h+1} \left( \Psi_{t,i}^{K_t,h}(M_{t-h:t}) + \widetilde{A}_{K_{t:t-i+1}} \mathbf{1}_{i=h+1} + \widetilde{A}_{K_{t:t-h}} B_{t-h-1} M_{t-h-1}^{[i-h-2]} \mathbf{1}_{1\leq i-h-1\leq H} \right) w_{t-i} \\
&= \widetilde{A}_{K_{t:t-h-1}} x_{t-h-1} + \sum_{i=0}^{H+h+1} \Psi_{t,i}^{K_t,h+1}(M_{t-h-1:t}) w_{t-i}.
\end{aligned}
\tag{99}
$$

$\square$

### D.5. Sufficiency of DAC Policy

**Proposition 2** (Sufficiency of DAC Policy). *Suppose the initial state is $x_0 = 0$, for a sequence of $(\kappa, \gamma)$ strongly stable time-varying controllers $K_0^*, \ldots, K_t^*$, there exist a policy $\pi_t(K_t, M_t^*)$, with $M_t^{*, [i]} \triangleq (K_t - K_t^*)\widetilde{A}_{K_{t-1:t-i}^*}$, where $\widetilde{A}_{K_{t:t-i}^*} \triangleq \prod_{\tau=t}^{t-i}(A_\tau - B_\tau K_\tau^*)$ and $\widetilde{A}_{K_{t:t-i}^*} \triangleq \mathbf{I}$ if $i < 0$, such that*

$$\sum_{t=0}^{T}\left(c_t\left(x_t\left(M_{0:t}^*\right), u_t\left(M_{0:t}^*\right)\right) - c_t\left(x_t^*\left(0\right), u_t^*\left(0\right)\right)\right) \leq \frac{4TG_cDWH\kappa_B^2\kappa^6(1-\gamma)^{H-1}}{\gamma}. \tag{100}$$

*Proof.* The proof is analogous to (Agarwal et al., 2019, Lemma 5.2), but adjusts the definition of $M_{0:t}^*$ to handle the time-varying controllers.

By definition, the state propagated by the sequence of time-varying controller $K_1^*, \ldots, K_t^*$ is

$$x_{t+1}^*(0) = \sum_{i=0}^{t} \widetilde{A}_{K_{t:t-i+1}^*} w_{t-i}. \tag{101}$$

Consider the state transition matrix $\Psi_{t,i}^{K_t,h}\left(M_{t-h:t}^*\right)$ for any $i \leq H$ and $H \leq h$. With eq. (97), we have

$$
\begin{aligned}
\Psi_{t,i}^{K_t,h}\left(M_{t-h:t}^*\right) &= \widetilde{A}_{K_{t:t-i+1}} + \sum_{j=1}^{i} \widetilde{A}_{K_{t:t-i+j+1}} B_{t-i+j}\left(K_{t-i+j} - K_{t-i+j}^*\right) \widetilde{A}_{K_{t-i+j-1:t-i+1}}^* \\
&= \widetilde{A}_{K_{t:t-i+1}} + \sum_{j=1}^{i} \widetilde{A}_{K_{t:t-i+j+1}}\left(\widetilde{A}_{K_{t-i+j}^*} - \widetilde{A}_{K_{t-i+j}}\right) \widetilde{A}_{K_{t-i+j-1:t-i+1}}^* \\
&= \widetilde{A}_{K_{t:t-i+1}} + \sum_{j=1}^{i}\left(\widetilde{A}_{K_{t:t-i+j+1}}\widetilde{A}_{K_{t-i+j:t-i+1}^*} - \widetilde{A}_{K_{t:t-i+j}}\widetilde{A}_{K_{t-i+j-1:t-i+1}}^*\right) \\
&= \widetilde{A}_{K_{t:t-i+1}^*},
\end{aligned}
\tag{102}
$$

where the last equality holds by telescoping the summation. Therefore, by setting $h = t$ in eq. (94), we have

$$x_{t+1}\left(M_{0:t}^*\right) = \sum_{i=0}^{H} \widetilde{A}_{K_{t:t-i+1}^*} w_{t-i} + \sum_{i=H+1}^{t} \Psi_{t,i}^{K_t,t}\left(M_{0:t}^*\right) w_{t-i}. \tag{103}$$

Combine eqs. (101) and (103), we get

$$
\begin{aligned}
\left\|x_{t+1}^*(0) - x_{t+1}\left(M_{0:t}^*\right)\right\|_2 &\leq W\left(\sum_{i=H+1}^{t}\left\|\Psi_{t,i}^{K_t,t}\left(M_{0:t}^*\right)\right\|_{\mathrm{op}} + \sum_{i=H+1}^{t}\left\|\widetilde{A}_{K_{t:t-i+1}^*}\right\|_{\mathrm{op}}\right) \\
&\leq W\left(\sum_{i=H+1}^{t}\left(2\kappa^2(1-\gamma)^i + H\kappa_B^2\kappa^5(1-\gamma)^{i-1}\right)\right) \\
&\leq W\left(2\kappa^2(1-\gamma)^{H+1}\gamma^{-1} + H\kappa_B^2\kappa^5(1-\gamma)^H\gamma^{-1}\right) \\
&\leq \kappa^2 W(1-\gamma)^{H+1}\gamma^{-1}\left(2(1-\gamma) + H\kappa_B^2\kappa^3\right) \\
&\leq H\kappa_B^2\kappa^5 W(1-\gamma)^H\gamma^{-1}\left(2(1-\gamma) + 1\right) \\
&\leq 2H\kappa_B^2\kappa^5 W(1-\gamma)^H\gamma^{-1},
\end{aligned}
\tag{104}
$$

where the second inequality holds due to Lemma 7.

Similarly, we analyze the difference in control input,

$$
\begin{aligned}
\left\| u_{t+1}^*(0) - u_{t+1}\left(M_{0:t+1}^*\right) \right\|_2 &= \left\| -K_{t+1}^* x_{t+1}^* - \left( -K_{t+1} x_{t+1}(M_{0:t}^*) + \sum_{i=1}^{H} M_{t+1}^{*,[i-1]} w_{t+1-i} \right) \right\|_2 \\
&= \left\| -K_{t+1}^* x_{t+1}^* + K_{t+1} x_{t+1}(M_{0:t}^*) - \sum_{i=1}^{H} (K_{t+1} - K_{t+1}^*) \widetilde{A}_{K_{t:t-i+2}^*} w_{t+1-i} \right\|_2 \\
&= \left\| -K_{t+1}^* \left( x_{t+1}^* - \sum_{i=0}^{H-1} \widetilde{A}_{K_{t:t-i+1}^*} w_{t-i} \right) + K_{t+1} \left( x_{t+1}(M_{0:t}^*) - \sum_{i=0}^{H-1} \widetilde{A}_{K_{t:t-i+1}^*} w_{t-i} \right) \right\|_2 \\
&= \left\| -K_{t+1}^* \sum_{i=H}^{t} \widetilde{A}_{K_{t:t-i+1}^*} w_{t-i} + K_{t+1} \sum_{i=H}^{t} \Psi_{t,i}^{K_t,t}(M_{0:t}^*) w_{t-i} \right\|_2 \\
&\leq 2H\kappa_B^2 \kappa^6 W (1-\gamma)^H \gamma^{-1}.
\end{aligned}
$$

(105)

Using eqs. (104) and (105), the Lipschitz assumption (Assumption 6), and the boundedness result (Lemma 8), we complete the proof. □

## D.6. Approximation of Truncated Loss

**Theorem 8** (Approximation of Truncated Loss: Theorem 5.3 of (Agarwal et al., 2019)). *Define* $D \triangleq \frac{W\kappa^3(1+H\kappa_B\tau)}{\gamma(1-\kappa^2(1-\gamma)^{H+1})} + \frac{W\tau}{\gamma}$. *For any* $(\kappa,\gamma)$ *strongly stable linear controller* $K_t$ *at iteration* $t$, *and any* $\tau > 0$ *such that the sequence of* $M_0, \ldots, M_T$ *satisfies* $\left\| M_t^{[i]} \right\|_{\text{op}} \leq \tau(1-\gamma)^i$, *the approximation error between original loss and truncated loss is at most*

$$
\left| \sum_{t=0}^{T} \left( c_t\left(x_t\left(M_{0:t-1}\right), u_t\left(M_{0:t}\right)\right) - f_t\left(M_{t-1-H:t}\right) \right) \right| \leq 2TG_c D^2 \kappa^3 (1-\gamma)^{H+1}.
$$

(106)

*Proof.* By the Lipschitz continuity and definition of the truncated loss, we get that

$$
\begin{aligned}
&c_t\left(x_t\left(M_{0:t-1}\right), u_t\left(M_{0:t}\right)\right) - f_t\left(M_{t-H-1:t}\right) \\
=&c_t\left(x_t\left(M_{0:t-1}\right), u_t\left(M_{0:t}\right)\right) - c_t\left(y_t\left(M_{t-H-1:t-1}\right), v_t\left(M_{t-H-1:t}\right)\right) \\
\leq&G_c D\left(\left\| x_t\left(M_{0:t-1}\right) - y_t\left(M_{t-H-1:t-1}\right)\right\|_2 + \left\| u_t\left(M_{0:t}\right) - v_t\left(M_{t-H-1:t}\right)\right\|_2\right) \\
\leq&G_c D\left(\kappa^2(1-\gamma)^{H+1}D + \kappa^3(1-\gamma)^{H+1}D\right) \\
\leq&2G_c D^2 \kappa^3 (1-\gamma)^{H+1},
\end{aligned}
$$

(107)

where the first inequality uses the Lipschitz assumption Assumption 6 and the second inequality uses boundedness in Lemma 8. The result in eq. (106) is obtained by summing over the iterations from $t = 1$ to $T$. □

## D.7. Supporting Lemmas

In this part, we provide supporting lemmas used in the analysis of online non-stochastic control. In particular,

- Lemma 6 presents the relationship between the $\ell_1$, op norm and Frobenius norm in the $\mathcal{M}$-space.

- Lemma 7 shows the norm of transfer matrix in eq. (95) is upper bounded.

- Lemma 8 provides the boundedness of several variables of interest.

- Lemma 9 shows properties of the truncated functions $\{f_t\}$ and the feasible set $\mathcal{M}$.

**Lemma 6** (Norm Relations over $\mathcal{M}$-space). *For any* $M = \left(M^{[0]}, \ldots, M^{[H-1]}\right) \in \mathcal{M} \subseteq \left(\mathbb{R}^{d_u \times d_x}\right)^H$, *its* $\ell_1$, *op norm and Frobenius norm are defined by*

$$
\|M\|_{\ell_1,\text{op}} \triangleq \sum_{i=0}^{H-1} \left\| M^{[i]} \right\|_{\text{op}}, \text{ and } \|M\|_{\text{F}} \triangleq \sqrt{\sum_{i=0}^{H-1} \left\| M^{[i]} \right\|_{\text{F}}^2}.
$$

(108)

*We then have the following inequalities on their relations:*

$$\|M\|_{\ell_1,\mathrm{op}} \leq \sqrt{H}\|M\|_{\mathrm{F}}, \text{ and } \|M\|_{\mathrm{F}} \leq \sqrt{d}\|M\|_{\ell_1,\mathrm{op}}, \tag{109}$$

*where* $d = \min\{d_u, d_x\}$.

*Proof.* For any matrix $X \in \mathbb{R}^{m \times n}$,

$$\|X\|_{\mathrm{op}} \leq \|X\|_{\mathrm{F}} \leq \sqrt{\min\{m,n\}}\|X\|_{\mathrm{op}}. \tag{110}$$

Therefore, by definition and Cauchy-Schwarz inequality, we obtain

$$\|M\|_{\ell_1,\mathrm{op}} = \sum_{i=0}^{H-1}\left\|M^{[i]}\right\|_{\mathrm{op}} \leq \sum_{i=0}^{H-1}\left\|M^{[i]}\right\|_{\mathrm{F}} \leq \sqrt{H}\|M\|_{\mathrm{F}}. \tag{111}$$

On the other hand, we have

$$\|M\|_{\mathrm{F}} = \sqrt{\sum_{i=0}^{H-1}\left\|M^{[i]}\right\|_{\mathrm{F}}^2} \leq \sum_{i=0}^{H-1}\left\|M^{[i]}\right\|_{\mathrm{F}} \leq \sum_{i=0}^{H-1}\sqrt{d}\left\|M^{[i]}\right\|_{\mathrm{op}} = \sqrt{d}\|M\|_{\ell_1,\mathrm{op}}. \tag{112}$$

$\square$

**Lemma 7** (Bounded Transfer Matrix). *Suppose $K_t$ is $(\kappa, \gamma)$-strongly stable at each iteration $t$. Suppose that for every $i \in \{0, \ldots, H-1\}$ and every $t \in \{1, \ldots, T\}$, we have $\left\|M_t^{[i]}\right\|_{\mathrm{op}} \leq \tau(1-\gamma)^i$ for some $\tau > 0$. Then, the transfer matrix is bounded as*

$$\left\|\Psi_{t,i}^{K,h}\right\|_{\mathrm{op}} \leq \kappa^2(1-\gamma)^i \mathbf{1}_{i\leq h} + H\kappa_B\kappa^2\tau(1-\gamma)^{i-1}. \tag{113}$$

*Proof.* We follow the proof of (Agarwal et al., 2019, Lemma 5.4). By definition of the transfer matrix $\Psi_{t,i}^{K,h}$ in eq. (95), we have

$$\left\|\Psi_{t,i}^{K,h}\right\|_{\mathrm{op}} = \left\|\widetilde{A}_{K_{t:t-i+1}}\mathbf{1}_{i\leq h} + \sum_{j=0}^{h}\widetilde{A}_{K_{t:t-j+1}}B_{t-j}M_{t-j}^{[i-j-1]}\mathbf{1}_{1\leq i-j\leq H}\right\|_{\mathrm{op}}$$

$$\leq \left\|\widetilde{A}_{K_{t:t-i+1}}\right\|_{\mathrm{op}}\mathbf{1}_{i\leq h} + \sum_{j=0}^{h}\left\|\widetilde{A}_{K_{t:t-j+1}}\right\|_{\mathrm{op}}\|B_{t-j}\|_{\mathrm{op}}\left\|M_{t-j}^{[i-j-1]}\right\|_{\mathrm{op}}\mathbf{1}_{1\leq i-j\leq H}$$

$$\leq \kappa^2(1-\gamma)^i\mathbf{1}_{i\leq h} + \sum_{j=0}^{H-1}\kappa^2(1-\gamma)^j\kappa_B\tau(1-\gamma)^{i-j-1} \tag{114}$$

$$\leq \kappa^2(1-\gamma)^i\mathbf{1}_{i\leq h} + \kappa^2\kappa_B\tau\sum_{j=1}^{H}(1-\gamma)^{i-1}$$

$$= \kappa^2(1-\gamma)^i\mathbf{1}_{i\leq h} + H\kappa^2\kappa_B\tau(1-\gamma)^{i-1}. \qquad \square$$

**Lemma 8** (Bounded State and Control). *Suppose $K_t$ and $K_t^*$ are $(\kappa, \gamma)$-strongly stable linear controllers at each iteration $t \in \{1, \ldots, T\}$. Suppose that for every $i \in \{0, \ldots, H-1\}$ and every $t \in \{1, \ldots, T\}$, we have $\left\|M_t^{[i]}\right\|_{\mathrm{op}} \leq \tau(1-\gamma)^i$ for some $\tau > 0$. Define $D \triangleq \frac{W(\kappa^3 + H\kappa_B\kappa^3\tau)}{\gamma(1-\kappa^2(1-\gamma)^{H+1})} + \frac{W\tau}{\gamma}$. Then, we have*

$$\|x_t(M_{0:t-1})\|_2 \leq D, \|y_t(M_{t-H-1:t-1})\|_2 \leq D, \left\|x_t^{K^*}\right\|_2 \leq D; \tag{115}$$

$$\|u_t(M_{0:t})\|_2 \leq D, \|v_t(M_{t-H-1:t})\|_2 \leq D; \tag{116}$$

$$\|x_t(M_{0:t-1}) - y_t(M_{t-1-H:t-1})\|_2 \leq \kappa^2(1-\gamma)^{H+1}D; \tag{117}$$

$$\|u_t(M_{0:t}) - v_t(M_{t-1-H:t})\|_2 \leq \kappa^3(1-\gamma)^{H+1}D. \tag{118}$$

*Proof.* The proof is analogous to (Agarwal et al., 2019, Lemma 5.5). We first consider eq. (115):

$$
\begin{aligned}
\|x_t(M_{0:t-1})\|_2 &= \left\| \widetilde{A}_{K_{t-1:t-H-1}} x_{t-H-1}(M_{0:t-H-2}) + \sum_{i=0}^{2H} \Psi_{t-1,i}^{K_t,H}(M_{t-H-1:t-1}) w_{t-i-1} \right\|_2 \\
&\le \kappa^2(1-\gamma)^{H+1} \|x_{t-H-1}(M_{0:t-H-2})\|_2 + W \sum_{i=0}^{2H} \left\| \Psi_{t-1,i}^{K_t,H}(M_{t-H-1:t-1}) \right\|_{\text{op}} \\
&\le \kappa^2(1-\gamma)^{H+1} \|x_{t-H-1}(M_{0:t-H-2})\|_2 + W \sum_{i=0}^{2H} \left( \kappa^2(1-\gamma)^i + H\kappa_B\kappa^2\tau(1-\gamma)^{i-1} \right) \\
&\le \kappa^2(1-\gamma)^{H+1} \|x_{t-H-1}(M_{0:t-H-2})\|_2 + W\frac{\kappa^2 + H\kappa_B\kappa^2\tau}{\gamma} \\
&\le \frac{W\left(\kappa^2 + H\kappa_B\kappa^2\tau\right)}{\gamma\left(1-\kappa^2(1-\gamma)^{H+1}\right)} \le D,
\end{aligned}
\tag{119}
$$

where the fourth inequality is a summation of geometric series and the ratio of this series is $\kappa^2(1-\gamma)^{H+1}$.

Similarly, we have

$$
\begin{aligned}
\|y_t(M_{t-H-1:t-1})\|_2 &= \left\| \sum_{i=0}^{2H} \Psi_{t-1,i}^{K_t,H}(M_{t-H-1:t-1}) w_{t-i-1} \right\|_2 \\
&\le W \sum_{i=0}^{2H} \left\| \Psi_{t-1,i}^{K_t,H}(M_{t-H-1:t-1}) \right\|_{\text{op}} \\
&\le W \sum_{i=0}^{2H} \left( \kappa^2(1-\gamma)^i + H\kappa_B\kappa^2\tau(1-\gamma)^{i-1} \right) \\
&\le W\frac{\kappa^2 + H\kappa_B\kappa^2\tau}{\gamma} \\
&\le \frac{W\left(\kappa^2 + H\kappa_B\kappa^2\tau\right)}{\gamma} \le D,
\end{aligned}
\tag{120}
$$

and

$$
\|x_t^*\|_2 = \left\| \sum_{i=0}^{t-1} \widetilde{A}_{K_{t-1:t-i}^*} w_{t-i-1} \right\|_2 \le W \sum_{i=0}^{t-1} \kappa^2(1-\gamma)^i \le \frac{W\kappa^2}{\gamma} \le D.
\tag{121}
$$

Next, we can show eq. (117) as follows,

$$
\|x_t(M_{0:t-1}) - y_t(M_{t-H-1:t-1})\|_2 = \left\| \widetilde{A}_{K_{t-1:t-H-1}} x_{t-H-1}(M_{0:t-H-2}) \right\|_2 \le \kappa^2(1-\gamma)^{H+1} D.
\tag{122}
$$

We now consider eqs. (116) and (118):

$$
\begin{aligned}
\|u_t(M_{0:t})\|_2 &= \left\| -K_t x_t(M_{0:t-1}) + \sum_{i=1}^{H} M_t^{[i-1]} w_{t-i} \right\|_2 \\
&\le \kappa \|x_t(M_{0:t-1})\|_2 + \sum_{i=1}^{H} W\tau(1-\gamma)^{i-1} \\
&\le \frac{W\left(\kappa^3 + H\kappa_B\kappa^3\tau\right)}{\gamma\left(1-\kappa^2(1-\gamma)^{H+1}\right)} + \frac{W\tau}{\gamma} \le D,
\end{aligned}
\tag{123}
$$

$$
\|v_t(M_{t-H-1:t})\|_2 \le \kappa \|y_t(M_{t-H-1:t-1})\|_2 + \sum_{i=1}^{H} W\tau(1-\gamma)^{i-1} \le D,
\tag{124}
$$

and

$$\left\| u_t \left( M_{0:t-1} \right) - v_t \left( M_{t-H-1:t-1} \right) \right\|_2 = \left\| -K_t \left( x_t \left( M_{0:t-1} \right) - y_t \left( M_{t-H-1:t-1} \right) \right) \right\|_2 \leq \kappa^3 (1-\gamma)^{H+1} D. \tag{125}$$

$\square$

**Lemma 9.** *Define* $D \triangleq \frac{W\kappa^3(1+H\kappa_B\tau)}{\gamma(1-\kappa^2(1-\gamma)^{H+1})} + \frac{W\tau}{\gamma}$. *The truncated loss* $f_t : \mathcal{M}^{H+2} \mapsto \mathbb{R}$ *is* $L_f$-*coordinate-wise Lipschitz and has bounded gradient norm* $G_f$. *In addition, the radius of feasible set in* $\mathcal{M}$-*space is bounded by* $D_f$. *Formally,*

1. *The truncated loss function is* $L_f$-*coordinate-wise Lipschitz with respect to the Euclidean (i.e., Frobenius) norm, i.e.,*

$$\left| f_t \left( M_{t-H-1}, \ldots, M_{t-k}, \ldots, M_t \right) - f_t \left( M_{t-H-1}, \ldots, \widetilde{M}_{t-k}, \ldots, M_t \right) \right| \leq L_f \left\| M_{t-k} - \widetilde{M}_{t-k} \right\|_{\mathrm{F}}, \tag{126}$$

*where* $L_f \leq 3 G_c D \sqrt{H} \kappa_B \kappa^3 (1-\gamma)^{k-1} W$.

2. *The gradient norm of surrogate loss* $\widetilde{f}_t : \mathcal{M} \mapsto \mathbb{R}$ *is bounded by* $G_f$, *i.e.,* $\left\| \nabla_M \widetilde{f}_t(M) \right\|_{\mathrm{F}} \leq G_f$ *holds for any* $M \in \mathcal{M}$ *and any* $t \in \{1, \ldots, T\}$, *where* $G_f \leq 3 H d_x d_u G_c W \kappa_B \kappa^3 \gamma^{-1}$.

3. *The diameter of the feasible set is at most* $D_f$, *i.e.,* $\|M - M'\|_{\mathrm{F}} \leq D_f$ *holds for any* $M, M' \in \mathcal{M}$, *where* $D_f \leq 2\sqrt{d} \kappa_B \kappa^3 \gamma^{-1}$.

*Proof.* The proof follows the steps in (Agarwal et al., 2019, Lemma 5.6 and 5.7) and (Zhao et al., 2022, Lemma 20). We first prove claim 1. To this end, we use the following notations:

$$\begin{aligned}
M_{t-H-1:t} &\triangleq \{ M_{t-H-1} \ldots M_{t-k} \ldots M_t \}; \\
M_{t-H-1:t-1} &\triangleq \{ M_{t-H-1} \ldots M_{t-k} \ldots M_{t-1} \}; \\
\widetilde{M}_{t-H-1:t} &\triangleq \left\{ M_{t-H-1} \ldots \widetilde{M}_{t-k} \ldots M_t \right\}; \\
\widetilde{M}_{t-H-1:t-1} &\triangleq \left\{ M_{t-H-1} \ldots \widetilde{M}_{t-k} \ldots M_{t-1} \right\}; \\
y_t &\triangleq y_t \left( M_{t-H-1:t-1} \right); \\
\widetilde{y}_t &\triangleq y_t \left( \widetilde{M}_{t-H-1:t-1} \right); \\
v_t &\triangleq v_t \left( M_{t-H-1:t} \right); \\
\widetilde{v}_t &\triangleq v_t \left( \widetilde{M}_{t-H-1:t} \right).
\end{aligned} \tag{127}$$

By definition of $f_t$, we have

$$f_t \left( M_{t-H-1:t} \right) - f_t \left( \widetilde{M}_{t-H-1:t} \right) = c_t \left( y_t, v_t \right) - c_t \left( \widetilde{y}_t, \widetilde{v}_t \right) \leq G_c D \left\| y_t - \widetilde{y}_t \right\|_2 + G_c D \left\| v_t - \widetilde{v}_t \right\|_2. \tag{128}$$

Then consider $\| y_t - \widetilde{y}_t \|$ and $\| v_t - \widetilde{v}_t \|$:

$$\begin{aligned}
\left\| y_t^K - \widetilde{y}_t^K \right\|_2 &= \left\| \sum_{i=0}^{2H} \left( \Psi_{t-1,i}^{K_t,H} \left( M_{t-H-1:t-1} \right) - \Psi_{t-1,i}^{K_t,H} \left( \widetilde{M}_{t-H-1:t-1} \right) \right) w_{t-1-i} \right\|_2 \\
&= \left\| \widetilde{A}_{K_{t-1:t-k+1}} B_{t-k} \sum_{i=0}^{2H} \left( M_{t-k}^{[i-k]} - \widetilde{M}_{t-k}^{[i-k]} \right) \mathbf{1}_{0 \leq i-k \leq H-1} w_{t-1-i} \right\|_2 \\
&\leq \kappa_B \kappa^2 (1-\gamma)^{k-1} W \sum_{i=1}^{H} \left\| M_{t-k}^{[i-1]} - \widetilde{M}_{t-k}^{[i-1]} \right\|_{\mathrm{op}} \\
&\leq \sqrt{H} \kappa_B \kappa^2 (1-\gamma)^{k-1} W \left\| M_{t-k} - \widetilde{M}_{t-k} \right\|_{\mathrm{F}},
\end{aligned} \tag{129}$$

and we have

$$
\begin{aligned}
\|v_t - \widetilde{v}_t\|_2 &= \left\| -K\left(y_t - \widetilde{y}_t\right) + \mathbf{1}_{\{k=0\}} \sum_{i=1}^{H} \left(M_{t-k}^{[i-1]} - \widetilde{M}_{t-k}^{[i-1]}\right) w_{t-i} \right\|_2 \\
&\le \left(\sqrt{H}\kappa_B \kappa^3 (1-\gamma)^{k-1} W + \sqrt{H} W\right) \left\| M_{t-k} - \widetilde{M}_{t-k} \right\|_{\mathrm{F}} \\
&\le 2\sqrt{H}\kappa_B \kappa^3 (1-\gamma)^{k-1} W \left\| M_{t-k} - \widetilde{M}_{t-k} \right\|_{\mathrm{F}}.
\end{aligned}
\tag{130}
$$

Combining the above equations, we obtain

$$
\begin{aligned}
f_t\left(M_{t-H-1:t}\right) - f_t\left(\widetilde{M}_{t-H-1:t}\right) &\le G_c D \left\| y_t^K - \widetilde{y}_t^K \right\|_2 + G_c D \left\| v_t^K - \widetilde{v}_t^K \right\|_2 \\
&\le G_c D \sqrt{H} \kappa_B \kappa^2 (1-\gamma)^{k-1} W \left\| M_{t-k} - \widetilde{M}_{t-k} \right\|_{\mathrm{F}} \\
&\quad + 2 G_c D \kappa_B \kappa^3 (1-\gamma)^{k-1} W \left\| M_{t-k} - \widetilde{M}_{t-k} \right\|_{\mathrm{F}} \\
&\le 3 G_c D \sqrt{H} \kappa_B \kappa^3 (1-\gamma)^{k-1} W \left\| M_{t-k} - \widetilde{M}_{t-k} \right\|_{\mathrm{F}}.
\end{aligned}
\tag{131}
$$

Therefore, we have $L_f \le 3 G_c D \sqrt{H} \kappa_B \kappa^3 (1-\gamma)^{k-1} W$.

Now consider claim 2. We need to bound $\nabla_{M_{p,q}^{[r]}} \widetilde{f}_t(M)$ for every $p \in \{1, \ldots, d_u\}$, $q \in \{1, \ldots, d_x\}$, and $r \in \{0, \ldots, H-1\}$,

$$
\left| \nabla_{M_{p,q}^{[r]}} \widetilde{f}_t(M) \right| \le G_c \left\| \frac{\partial y_t(M)}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}} + G_c \left\| \frac{\partial v_t(M)}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}}.
\tag{132}
$$

Now we aim to bound the two terms of the right-hand side respectively:

$$
\begin{aligned}
\left\| \frac{\partial y_t(M)}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}} &\le \left\| \sum_{i=0}^{2H} \sum_{j=0}^{H} \left[ \frac{\partial \widetilde{A}_{K_{t:t-j+1}} B_{t-j} M^{[i-j-1]}}{\partial M_{p,q}^{[r]}} \right] w_{t-1-i} \mathbf{1}_{0 \le i-j \le H-1} \right\|_{\mathrm{F}} \\
&\le W \kappa_B \kappa^2 \left\| \frac{\partial M^{[r]}}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}} \sum_{i=r+1}^{r+H+1} (1-\gamma)^{i-r-1} \\
&\le \frac{W \kappa_B \kappa^2}{\gamma} \left\| \frac{\partial M^{[r]}}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}} \\
&\le \frac{W \kappa_B \kappa^2}{\gamma}; \\
\left\| \frac{\partial v_t(M)}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}} &\le \kappa \left\| \frac{\partial y_t(M)}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}} + \sum_{i=1}^{H} \left\| \frac{\partial M^{[i-1]}}{\partial M_{p,q}^{[r]}} w_{t-i} \right\|_{\mathrm{F}} \\
&\le \frac{W \kappa_B \kappa^3}{\gamma} + W \left\| \frac{\partial M^{[r]}}{\partial M_{p,q}^{[r]}} \right\|_{\mathrm{F}} \\
&\le W \left( \frac{\kappa_B \kappa^3}{\gamma} + 1 \right).
\end{aligned}
\tag{133}
$$

Therefore, we have

$$
\left| \nabla_{M_{p,q}^{[r]}} \widetilde{f}_t(M) \right| \le G_c \frac{W \kappa_B \kappa^2}{\gamma} + G_c W \left( \frac{\kappa_B \kappa^3}{\gamma} + 1 \right) \le 3 G_c W \kappa_B \kappa^3 \gamma^{-1}.
\tag{134}
$$

Thus, $\left\| \nabla_M \widetilde{f}_t(M) \right\|_{\mathrm{F}}$ is at most $3 H d_x d_u G_c W \kappa_B \kappa^3 \gamma^{-1}$.

*Table 2.* Comparison of the **OGD** (Zinkevich, 2003), **Ader** (Zhang et al., 2018), **Scream** (Zhao et al., 2022), and **Meta-OFW** algorithms in terms of cumulative loss for 10000 time steps. The blue numbers correspond to the best performance and the red numbers correspond to the worse.

| Noise Distribution | Sinusoidal Weights (eq. (137)) | | | | Step Weights (eq. (138)) | | | |
|---|---|---|---|---|---|---|---|---|
| | **Meta-OFW** | **Scream** | **Ader** | **OGD** | **Meta-OFW** | **Scream** | **Ader** | **OGD** |
| Gaussian | 15625 | 19725 | 21052 | 33574 | 9496 | 10704 | 11453 | 26790 |
| Uniform | 18299 | 93987 | 107096 | 30419 | 13395 | 39057 | 35313 | 39885 |
| Gamma | 16239 | 16138 | 18039 | 17484 | 9184 | 61989 | 75505 | 45398 |
| Beta | 21448 | 34146 | 30990 | 30253 | 15982 | 29301 | 30799 | 28859 |
| Exponential | 10621 | 254815 | 252227 | 28859 | 4366 | 227860 | 204844 | 53626 |
| Weibull | 14068 | 91474 | 94040 | 38549 | 5623 | 182887 | 993734 | 92341 |

*Table 3.* Comparison of the **Scream** (Zhao et al., 2022) and **Meta-OFW** algorithms in terms of computational time and cumulative loss over 200 time steps for the case of Gaussian noise in eq. (136), sinusoidal weights in eq. (137), and across varying system dimensions $d_x$ and $d_u$.

| $(d_x,\ d_u)$ | Time (seconds) | | Cumulative Loss | |
|---|---|---|---|---|
| | **Meta-OFW** | **Scream** | **Meta-OFW** | **Scream** |
| (2, 1) | 52.49 | 17.64 | 915.52 | 1819.41 |
| (4, 2) | 135.04 | 177.47 | 1388.56 | 3231.89 |
| (6, 3) | 287.67 | 605.48 | 1235.83 | 2021.47 |
| (8, 4) | 504.82 | 1351.29 | 1421.05 | 1873.79 |
| (10, 5) | 786.87 | 2219.85 | 1202.43 | 1439.22 |
| (12, 6) | 998.22 | 3029.89 | 893.17 | 984.96 |
| (14, 7) | 2022.5 | 5531.36 | 797.80 | 958.09 |

Finally, we prove claim 3. By construction of $M^{[i]}, \forall i \in \{0, \ldots, H-1\}$, we require $\|M^{[i]}\|_{\mathrm{op}} \leq \kappa_B \kappa^3 (1-\gamma)^i$. Therefore, utilizing Lemma 6 we have

$$
\begin{aligned}
\max_{M_1,M_2\in\mathcal{M}} \|M_1 - M_2\|_{\mathrm{F}} &\leq \sqrt{d} \max_{M_1,M_2\in\mathcal{M}} \|M_1 - M_2\|_{\ell_1,\mathrm{op}} \\
&\leq \sqrt{d} \max_{M_1,M_2\in\mathcal{M}} \left( \|M_1\|_{\ell_1,\mathrm{op}} + \|M_2\|_{\ell_1,\mathrm{op}} \right) \\
&= \sqrt{d} \max_{M_1,M_2\in\mathcal{M}} \left( \sum_{i=0}^{H-1} \left\|M_1^{[i]}\right\|_{\mathrm{op}} + \left\|M_2^{[i]}\right\|_{\mathrm{op}} \right) \\
&\leq \sqrt{d} \max_{M_1,M_2\in\mathcal{M}} \left( 2\sum_{i=0}^{H-1} \kappa_B \kappa^3 (1-\gamma)^i \right) \\
&= 2\sqrt{d}\kappa_B \kappa^3 \sum_{i=0}^{H-1} (1-\gamma)^i \\
&\leq 2\sqrt{d}\kappa_B \kappa^3 \gamma^{-1}.
\end{aligned}
\tag{135}
$$

Hence, we finish the proof of all three claims in the statement. □

## E. Numerical Evaluations

We evaluate **Meta-OFW** (Algorithm 3) in simulated scenarios of online control of linear time-invariant systems.

**Compared Algorithms.** We compare **Meta-OFW** with the **OGD** (Zinkevich, 2003), **Ader** (Zhang et al., 2018), and **Scream** (Zhao et al., 2022) algorithms. All algorithms rely on the DAC policy (Agarwal et al., 2019).

**Simulation Setup.** We follow the setup as (Zhao et al., 2021) and consider linear systems of the form

$$
\begin{aligned}
x_{t+1} &= Ax_t + Bu_t + w_t \\
&= Ax_t + Bu_t + (\Delta_{t,A} x_t + \Delta_{t,B} u_t + \tilde{w}_t),
\end{aligned}
\tag{136}
$$

where $\tilde{w}_t$ and the elements of $\Delta_{t,A}$ and $\Delta_{t,B}$ are sampled from various distributions, specifically, Gaussian, Uniform, Gamma, Beta, Exponential, and Weibull distributions. We use memory length $H = 10$. The loss function has the form $c_t(x_t, u_t) = q_t x_t^\top x_t + r_t u_t^\top u_t$, where $q_t \in \mathbb{R}$ and $r_t \in \mathbb{R}$ are time-varying weights. Particularly, we consider two cases:

1. Sinusoidal weights defined as
$$
q_t = \sin(t/10\pi),\ r_t = \sin(t/20\pi).
\tag{137}
$$

2. Step weights defined as
$$
(q_t, r_t) =
\begin{cases}
\left(\frac{\log(2)}{2}, 1\right), & t \leq T/5, \\
(1, 1), & T/5 < t \leq 2T/5, \\
\left(\frac{\log(2)}{2}, \frac{\log(2)}{2}\right), & 2T/5 < t \leq 3T/5, \\
\left(1, \frac{\log(2)}{2}\right), & 3T/5 < t \leq 4T/5, \\
\left(\frac{\log(2)}{2}, 1\right), & 4T/5 < t \leq T.
\end{cases}
\tag{138}
$$

**Results.** We first compare **Meta-OFW** with the **OGD**, **Ader**, and **Scream** algorithms in terms of cumulative loss. The results are summarized in Table 2, showing that **Meta-OFW** achieved the lowest cumulative loss across all tested cases, except under gamma distribution with sinusoidal weights; in the best-case —exponential distribution with step weights— **Meta-OFW** is 52 times better than **Scream**.

We also vary the dimensions of the state $x_t$ and input the $u_t$, and compare **Meta-OFW** and **Scream** in terms of cumulative loss and computation time. The results are summarized in Table 3. As $d_x$ and $d_u$ increase, **Meta-OFW** is computationally three times faster than **Scream**, achieving also lower cumulative loss than **Scream** in all cases.