

---

# Delay-Adapted Policy Optimization and Improved Regret for Adversarial MDP with Delayed Bandit Feedback

---

Tal Lancewicki<sup>1</sup> Aviv Rosenberg<sup>2</sup> Dmitry Sotnikov<sup>2</sup>

## Abstract

Policy Optimization (PO) is one of the most popular methods in Reinforcement Learning (RL). Thus, theoretical guarantees for PO algorithms have become especially important to the RL community. In this paper, we study PO in adversarial MDPs with a challenge that arises in almost every real-world application – *delayed bandit feedback*. We give the first near-optimal regret bounds for PO in tabular MDPs, and may even surpass state-of-the-art (which uses less efficient methods). Our novel Delay-Adapted PO (DAPO) is easy to implement and to generalize, allowing us to extend our algorithm to: (i) infinite state space under the assumption of linear  $Q$ -function, proving the first regret bounds for delayed feedback with function approximation. (ii) deep RL, demonstrating its effectiveness in experiments on MuJoCo domains.

## 1. Introduction

Policy Optimization (PO) is one of the most widely-used methods in Reinforcement Learning (RL). It has demonstrated impressive empirical success (Levine & Koltun, 2013; Schulman et al., 2017; Haarnoja et al., 2018), leading to increasing interest in understanding its theoretical guarantees. While in recent years we have seen great advancement in theory of PO (Shani et al., 2020b; Luo et al., 2021; Chen et al., 2022b), our understanding is still very limited when considering *delayed feedback* – an important challenge that occurs in most practical applications. For example, recommendation systems often learn the utility of a recommendation based on the number of user conversions, which may happen with a variable delay after the recommendation was issued. Other notable examples in-

---

<sup>1</sup>Tel Aviv University (Research conducted while the author was an intern at Amazon Science) <sup>2</sup>Amazon Science. Correspondence to: Tal Lancewicki <lancewicki@mail.tau.ac.il>, Aviv Rosenberg <avivros007@gmail.com>.

clude communication between agents (Chen et al., 2020a), video streaming (Changuel et al., 2012) and robotics (Mahmood et al., 2018). To mitigate the gap in the PO literature, we study PO in the challenging adversarial MDP model (i.e., costs change arbitrarily) under bandit feedback with arbitrary unrestricted delays.

PO with delays was previously studied by Lancewicki et al. (2022b), but their regret bounds are far from optimal and scale with  $(K + D)^{2/3}$ , where  $K$  is the number of episodes and  $D$  is the total delay. Recently, Jin et al. (2022) achieved near-optimal  $\tilde{O}(H^2 S \sqrt{AK} + (HSA)^{1/4} H \sqrt{D})$  regret (ignoring logarithmic factors), where  $S$ ,  $A$  and  $H$  are the number of states, actions, and the episode length, respectively. However, their algorithm is not based on PO, but on the O-REPS method (Zimin & Neu, 2013) which requires solving a computationally expensive global optimization problem and cannot be extended to function approximation (FA). On the other hand, PO algorithms build on highly efficient local-search and extend naturally to FA (Tomar et al., 2022).

**Our Contributions.** In this paper, we vastly expand our understanding of PO and delayed feedback. We propose a novel Delay-Adapted PO method, called DAPO, which measures changes in the agent’s policy over the time of the delays and adapts its updates accordingly. First, we establish the power of DAPO in tabular MDPs, i.e., finite number of states and actions. We prove DAPO attains the first near-optimal regret bound for PO with delayed feedback  $\tilde{O}(H^3 S \sqrt{AK} + H^3 \sqrt{D})$ . This bound is tighter than (Jin et al., 2022) when the delay term is dominant and the number of states is significantly larger than the horizon (which is the common case). Moreover, it matches the lower bound of Lancewicki et al. (2022b) in the delay term up to factors of  $H$ , showing for the first time that the delay term in the regret does not need to scale with  $S$  or  $A$ . Importantly, if there is no delay, it matches the best known regret for PO (Luo et al., 2021).

Next, we show DAPO is easy to implement and naturally extends to function approximation in two important settings:

1. *Linear- $Q$* . We extend DAPO to MDPs with linear FA under standard assumptions (Luo et al., 2021) that  $Q$ -functions are linear in some known low-dimensional

features and also a simulator is available. We prove that DAPO achieves the first sub-linear regret for delayed feedback with FA, i.e., non-tabular MDP.

2. Deep RL. We show that the famous PPO algorithm (Schulman et al., 2017) can be easily combined with DAPO, and demonstrate superior empirical performance even in the presence of simple delays in experiments on MuJoCo domains (Todorov et al., 2012).

Throughout the paper we handle several technical challenges which are unique to PO algorithms with delayed feedback. The main challenge is to control the stability of the algorithm. This problem is more challenging in MDPs compared to multi-armed bandit (where there is no transition function to estimate), and is enhanced even further in Policy Optimization algorithms due to their local-search nature, as opposed to the global update of O-REPS methods – see more details in the proof sketch of Theorem 3.1. Further, the linear Q setting with delayed feedback was not studied before and requires a careful new algorithmic design and analysis. In particular, a-priori, it is highly unclear how to design a delay-adapted estimator and delay-adapted bonus term which sufficiently stabilize the algorithm – see more details in Section 4.

While the main contribution of this paper is the novel Delay-Adapted PO method, we also make substantial technical contributions that might be of independent interest. Our algorithms are based on PO with dilated bonuses (Luo et al., 2021), which dilate towards further horizons and do not satisfy standard Bellman equations. However, we are able to achieve the same regret guarantees without using dilated bonuses. Instead, we compute local bonuses and use them to construct a Q-function that operates as exploration bonuses. This has an important practical benefit – now bonuses can be approximated similarly to the Q-function. It also has a theoretical benefit – it greatly simplifies the analysis, making it more natural and easy to extend to new scenarios. Moreover, utilizing our new simplified analysis, we are able to give regret guarantees with high probability in the Linear setting and not just in expectation (as in Luo et al. (2021)). Finally, we also develop new analyses for handling delayed feedback when losses can be negative. This was not addressed in the delayed multi-armed bandit literature (or in previous papers on delays in MDPs), but cannot be avoided in our case since exploration bonuses are crucial to guaranteeing near-optimal regret but they might turn losses to negative.

### 1.1. Additional Related Work

Due to lack of space, this section only gives a brief overview of related work - for a full literature review see Appendix A. There is a rich literature on regret minimization in tabular MDPs, initiated with the seminal UCRL algorithm (Jaksch

et al., 2010) for stochastic losses that is based on the fundamental concept of Optimism Under Uncertainty. Their model was later extended to the more general adversarial MDP, where most algorithms are based on either the framework of occupancy measures (a.k.a, O-REPS) (Zimin & Neu, 2013; Jin et al., 2020a) or on the more practical PO (Even-Dar et al., 2009; Shani et al., 2020b). In recent years this line of research was extended beyond the tabular model to linear function approximation. For stochastic losses, existing algorithms are mostly based on optimism (Jin et al., 2020b), whereas most algorithms for adversarial losses are based on PO (Cai et al., 2020; Luo et al., 2021; Neu & Oikhovskaya, 2021) which extends much more naturally than O-REPS to function approximation. On the practical side, some of the most successful deep RL algorithms are built upon PO principles. These include the famous Trust Region PO (TRPO; Schulman et al. (2015)) as well as Proximal PO (PPO; Schulman et al. (2017)) which we will further discuss and adapt to delayed feedback in Section 5.

Regret minimization with delayed feedback was initially studied in Online Optimization and Multi-armed bandit (MAB) in both the stochastic setting (Agarwal & Duchi, 2012; Pike-Burke et al., 2018) and the adversarial setting (Cesa-Bianchi et al., 2016; Thune et al., 2019). As a natural extension, this line of work was generalized to delayed feedback in MDPs, where (Howson et al., 2021) consider the more restrictive stochastic model. Most related to our work are the works of Lancewicki et al. (2022b); Jin et al. (2022) that were mentioned earlier, and the work of Dai et al. (2022). They recently showed that Follow-The-Perturbed-Leader (FTPL) algorithms can also handle delayed feedback in adversarial MDPs. The efficiency of FTPL is similar to PO, but their regret bound is slightly weaker than Jin et al. (2022). Finally, a different line of work (Katsikopoulos & Engelbrecht, 2003; Walsh et al., 2009) consider delays in observing the current state. That setting is inherently different than ours (see Appendix A for more details).

## 2. Preliminaries

A finite-horizon episodic adversarial MDP is defined by a tuple  $M = (S; A; H; p; f; c; g_{k=1}^K)$ , where  $S$  and  $A$  are state and action spaces of sizes  $|S| = S$  and  $|A| = A$ , respectively,  $H$  is the horizon and  $K$  is the number of episodes.  $p : S \times A \rightarrow [0, 1]^S$  is the transition functions such that  $p_h(s'; s; a)$  is the probability to move to  $s'$  when taking action  $a$  in states at time  $h$ .  $f : S \times A \rightarrow [0, 1]^S$  are cost functions chosen by an oblivious adversary, where  $c_h^k(s; a)$  is the cost for taking action  $a$  at  $(s; h)$  in episode  $k$ .

A policy  $\pi : S \times [H] \rightarrow A$  is a function that gives the probability  $\pi_h(a|s)$  to take action  $a$  when visiting state  $s$  at time  $h$ . The value  $V_h(s; c)$  is the expected cost of  $\pi$  with respect to cost function  $c$  starting from  $s$  in time  $h$ , i.e.,

$V_h(s; c) = \mathbb{E} \sum_{t=0}^H c_{h^0}(s_{h^0}; a_{h^0}) \mid s_h = s$ , where the expectation is with respect to policy and transition function  $p$ , that is,  $a_{h^0} \sim p_{h^0}(\cdot \mid s_{h^0})$  and  $s_{h^0+1} \sim p_{h^0}(\cdot \mid s_{h^0}; a_{h^0})$ . The Q-function is defined by  $Q_h(s; a; c) = c_h(s; a) + \mathbb{E} p_h(\cdot \mid s; a; V_{h+1}(\cdot; c))$ , where  $\cdot$  is the dot product.

The learner interacts with the environment for  $K$  episodes. At the beginning of episode  $k$ , it picks a policy  $\pi^k$ , and starts in an initial state  $s_1^k = s_{\text{init}}$ . In each time step  $t \in [H]$ , it observes the current state  $s_t^k$ , draws an action from the policy  $a_t^k \sim \pi^k(\cdot \mid s_t^k)$  and transitions to the next state  $s_{t+1}^k \sim p_h(\cdot \mid s_t^k; a_t^k)$ . The feedback of episode  $k$  contains the cost function over the agent's trajectory  $c_h^k(s_h^k; a_h^k) g_{h=1}^H$ , i.e., bandit feedback. This feedback is observed only at the end of episode  $k + d^k$ , where the delays  $d^k g_{k=1}^K$  are unknown and chosen by the adversary together with the costs.

The goal of the learner is to minimize the regret, defined as the difference between the learner's cumulative expected cost and the best fixed policy in hindsight:

$$R_K = \sum_{k=1}^K V_1^k(s_{\text{init}}; c^k) - \min_{k=1}^K V_1(s_{\text{init}}; c^k);$$

In Section 3 we consider tabular MDPs i.e., MDPs with a finite number of states and actions. In Section 4 we consider the more general case that allows for infinite number of states but under the assumption that the Q-function is linear for all policies. We follow the standard definition (Abbasi-Yadkori et al., 2019; Neu & Olkhovskaya, 2021; Luo et al., 2021), which also assumes the number of actions is finite.

**Assumption 2.1 (Linear-Q).** Let  $\phi : S \times A \rightarrow [H] \times \mathbb{R}^n$  be a known feature mapping. Assume that for every episode  $k$ , policy  $\pi^k$  and step  $t$  there exist an unknown vector  $q_h^k \in \mathbb{R}^n$  such that  $Q_h(s; a; c^k) = \langle \phi_h(s; a), q_h^k \rangle$  for all  $(s; a)$ . Moreover,  $\|q_h(s; a)\|_2 \leq 1$  and  $\|q_h^k\|_2 \leq H \bar{n}$ .

**Additional notations.** Episode indices appear as superscripts and in-episode steps as subscripts. The total delay is  $D = \sum_k d^k$ , the maximal delay is  $d_{\max}$  and the number of episodes that their feedback arrives in the end of episode  $k$  is  $m^k = \#\{j : j + d^j = k\}$ . The occupancy measure  $q_h(s; a) = \Pr[s_h = s; a_h = a \mid s_1 = s_{\text{init}}]$  is the distribution that policy  $\pi$  induces over state-action pairs in step  $h$ , and  $q_h(s) = \sum_{a \in A} q_h(s; a)$ . The notations  $\mathcal{O}(\cdot)$  and  $\mathcal{H}(\cdot)$  hide poly-logarithmic factors including  $\log(K = \sum_k m^k)$  for confidence parameter,  $[n] = \{1; 2; \dots; n\}$  and the indicator of event  $E$  is  $\mathbb{1}\{E\}$ . Finally, denote by  $\pi^*$  the best fixed policy in hindsight and use the notation  $V_h^*(s)$ ,  $Q_h^*(s; a)$ ,  $c_h^*(s; a)$  when the policy and cost are  $\pi^*$  and  $c^*$ , respectively.

**Simplifying assumptions.** Similarly to Jin et al. (2022), we assume that  $K$ ,  $D$  and  $d_{\max}$  are known. This assumption

<sup>1</sup>If  $d^k = 0$ , we get standard online learning in adversarial MDP.

---

**Algorithm 1** DAPO with Known Transitions (Tabular)

---

Initialization: Set  $Q_h^1(a \mid s) = 1 - \epsilon$  for every  $(s; a; h)$ .  
 for  $k = 1; 2; \dots; K$  do  
   Play episode  $k$  with policy  $\pi^k$ .  
   # Policy Evaluation  
   for  $j$  such that  $j + d^j = k$  do  
     Observe bandit feedback  $c_h^j(s_h^j; a_h^j) g_{h=1}^H$ .  
     Compute delay-adapted estimator  $\hat{Q}_h^j(s; a)$  defined in Eq. (3).  
     Compute delay-adapted bonus  $B_h^j(s; a)$  as the Q-function of  $j$  with respect to the costs  $b_h^j(s)$  defined in Eq. (4). I.e., compute recursively for  $h = H; \dots; 1$ ,  $B_h^j(s; a) = b_h^j(s) + \mathbb{E}_{s^0 \sim p_h(\cdot \mid s; a); a^0 \sim p_{h+1}(\cdot \mid s^0)} B_{h+1}^j(s^0; a^0)$ .  
   end for  
   # Policy Improvement  
   Define the policy  $\pi^{k+1}$  for every  $(s; a; h)$  by:  
     
$$\pi^{k+1}(a \mid s) \propto \sum_{j: j + d^j = k} \exp(\hat{Q}_h^j(s; a) - B_h^j(s; a)) \quad (1)$$
  
   end for

---

simplifies presentation and can be easily removed without affecting the analysis. Bounding for delayed feedback (Bistriz et al., 2021; Lancewicki et al., 2022b). Bounds in the main text hide low-order terms and additive dependence on  $d_{\max}$  (see Remark C.2 on removing  $d_{\max}$  dependence). For full bounds see Appendix.

### 3. DAPO for Tabular MDP

In this section we present our novel Delayed-Adapted Policy Optimization algorithm (DAPO; presented in Algorithm 1) for the tabular case, where the number of states is finite. We use this fundamental model to develop a generic method for handling delayed feedback with Policy Optimization. Our approach consists of two important algorithmic features: a new delay-adapted importance-sampling estimator for the Q-function and a novel delay-adapted bonus term to drive exploration. Remarkably, this method extends naturally to both linear function approximation as we show in Section 4, and to deep RL with the extremely practical PPO algorithm (Schulman et al., 2017) as we show in Section 5.

To simplify presentation and focus on the contributions of our delay adaptation method, in this section we assume that the agent knows the transition function in advance. Generalizing DAPO to unknown transitions in the tabular case is fairly straightforward, and follows the common approach of optimism and confidence sets (Jaksch et al., 2010). Due to lack of space, in the main text we only provide sketches for the algorithms and proofs. The full versions (for both known and unknown transitions), together with the detailed analyses, can be found in Appendices B and C.

PO algorithms follow the algorithmic paradigm of Policy Iteration (see, e.g., Sutton & Barto (2018)). That is, in every iteration they perform an evaluation of the current policy and then a step of policy improvement. The improvement step is regularized to be “soft”, and is practically implemented by running an online multi-armed bandit algorithm, such as Hedge (Freund & Schapire, 1997), locally in each state. The losses that are fed to the algorithm are the estimated Q-functions, but in order to achieve the optimal regret, they are also combined with a bonus term which aims to stabilize the algorithm and drive exploration (keeping the estimated Q-function optimistic). The actual policy update step is based on exponential weights and presented in Eq. (1).

$$b_h^k(s) = \sum_{a \in \mathcal{A}} r_h^k(s; a) \frac{3 H \frac{k+d^k}{h} (a|s)}{c_h^k(s) \frac{k}{h} (a|s) + \epsilon}. \quad (4)$$

DAPO adapts to delays through the policy evaluation step. Remarkably, it adapts to delays near-optimally by computing the following simple ratio, which measures the local change in the agent’s policy through the time of the delay,

$$r_h^k(s; a) = \frac{\frac{k}{h} (a|s)}{\max_{a'} \frac{k}{h} (a'|s); \frac{k+d^k}{h} (a|s)}: \quad (2)$$

In order to get our delay-adapted Q-function estimation, we simply multiply  $r_h^k(s; a)$  by the standard importance-sampling estimator from the non-delayed setting (Luo et al., 2021). The result is:

$$\hat{Q}_h^k(s; a) = r_h^k(s; a) \frac{1}{c_h^k(s)} \frac{L_h^k(s; a)}{\frac{k}{h} (a|s)}; \quad (3)$$

where  $L_h^k = \sum_{h=0}^H c_{h0}^k(s_h^k; a_{h0}^k)$  is the realized cost-to-go from step  $h$  and  $\epsilon$  is an exploration parameter (Neu, 2015) needed to guarantee regret with high probability.

Intuitively, incorporating  $r_h^k(s; a)$  helps us control the variance of the estimator  $\hat{Q}_h^k(s; a)$  in the presence of delays, since it will be used only in episode  $k+d^k$  where actions are chosen according to  $\frac{k+d^k}{h}$  and not  $\frac{k}{h}$ . The maximum in the denominator is needed in order to keep the bias small. We note that this estimator is inspired by Jin et al. (2022), but there are two major differences. First, as Jin et al. (2022) perform their update globally in the space of state-action occupancy measures, their adaptation occurs in the state space as well. On the other hand we perform the update locally in each state, so our adaptation takes place only in the action space. Second, they directly change importance-sampling weighting, while we simply multiply standard estimators by the delay-adapted ratio. This seemingly minor nuance is critical in more complex (non-tabular) regimes, where its generality allows to utilize existing procedures from the non-delayed case (see more details in Sections 4 and 5).

Finally, to complement our new estimator, we devise an appropriate delay-adapted bonus  $b_h^k(s; a)$  based on the following delay-adapted local bonus (again obtained by combining  $r_h^k(s; a)$  with the original local bonus of Luo et al.

At this point, Luo et al. (2021) compute  $b_h^k(s; a)$  using a delayed Bellman equation that is not very intuitive. Instead, we compute  $b_h^k(s; a)$  with the regular Bellman equations, making it a proper Q-function. This is an important contribution that might be of independent interest for two reasons – theoretically the analysis becomes much simpler, and practically the bonuses can be approximated like a function.

Next, we present the regret guarantees of DAPO in tabular MDPs, and the key steps in the analysis, which highlight the intuition behind our algorithm design.

**Theorem 3.1.** Running DAPO in a tabular adversarial MDP guarantees with probability  $1 - \delta$ , for known transition,

$$R_K = \mathcal{O}(H^2 \overline{SAK} + H^3 \overline{K + D})$$

when setting  $\epsilon = H^2 \overline{SAK} + H^4(K + D)^{1=2}$  and  $\epsilon = 2 H$ , and for unknown transition,

$$R_K = \mathcal{O}(H^3 \overline{SAK} + H^3 \overline{D})$$

when setting  $\epsilon = H^2 \overline{SAK} + H^4(K + D)^{1=2}$  and  $\epsilon = 2 H$ .

This is a big improvement compared to the best known regret for PO (Lancewicki et al., 2022b) that scales as  $(K + D)^{2=3}$  (ignoring dependencies in  $H; S; A$ ). It is also better than the current state-of-the-art regret bound of Jin et al. (2022) in the case that there is significant delay and  $H$  (which occurs in almost every practical application). While their bound has better dependency in  $H$  (this is a known weakness of PO (Chen et al., 2022b)), we improve the dependency in  $S$  and  $A$ . Lancewicki et al. (2022b) also show a lower bound of  $(H^{3=2} \overline{SAK} + H \overline{D})$ . Thus, our bound shows for the first time that under the optimal regret, the delay term does not scale with  $S$  or  $A$ . The first term in our regret bound matches the state-of-the-art regret of non-delayed PO (Luo et al., 2021), and matches the best known regret for (non-delayed) adversarial MDPs in general up to factors of  $H$  (Jin et al., 2020a). Moreover, DAPO is the first efficient algorithm to be consistent with the optimal regret in delayed MAB, i.e., for  $H = S = 1$  we get the optimal regret of Thune et al. (2019); Bistritz et al. (2019). Finally, it is important to emphasize that PO algorithms are much more practical than O-REPS algorithms, and extend naturally to function approximation, as we show in Sections 4 and 5.

**Proof sketch of Theorem 3.1.** Much of the intuition for PO algorithms stems from a classic regret decomposition known as the value difference lemma (Even-Dar et al., 2009):

$R_K = \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^k(j, s) - h_h^k(j, s); Q_h^k(s); i$ . Fixing a state and step, the sum over  $k$  can be viewed as the regret of an online experts algorithm (e.g., Hedge) with respect to the losses  $Q_h^k(s); i$ . We propose to further decompose the regret as follows,

$$\begin{aligned}
 R_K &= \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^{k+d^k}(j, s); Q_h^k(s); i - \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^k(j, s); Q_h^k(s); i \\
 &+ \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^k(j, s); Q_h^k(s); i - \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^k(j, s); B_h^k(s); i \\
 &+ \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^{k+d^k}(j, s); Q_h^k(s); i - \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^{k+d^k}(j, s); B_h^k(s); i \\
 &+ \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^k(j, s); Q_h^k(s); i - \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^{k+d^k}(j, s); Q_h^k(s); i \\
 &+ \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^k(j, s); Q_h^k(s); i - \sum_{k,h} \mathbb{E}_s \sum_{q_h} h_h^k(j, s); B_h^k(s); i
 \end{aligned} \tag{5}$$

Indeed, the policy update step in Eq. (1) is an Hedge-style exponential weights update. This allows us to bound the term as it represents the regret of Hedge with respect to the losses  $Q_h^k(s); a - B_h^k(s); a$ . Note that the delayed feedback causes a shift of the agent's policies from  $h_h^k$  to  $h_h^{k+d^k}$ . As a result, we can bound (using Corollary E.7 in Appendix E):

$$\begin{aligned}
 \text{REG} &\leq \frac{H}{\epsilon} + \sum_{k,h;s;a} \mathbb{E}_s \sum_{q_h} \frac{h_h^{k+d^k}(a, j, s) (Q_h^k(s); a - B_h^k(s); a)^2}{h_h^{k+d^k}(a, j, s)} \\
 &\leq \frac{H}{\epsilon} + H^5 K + \sum_{k,h;s;a} \mathbb{E}_s \sum_{q_h} \frac{h_h^{k+d^k}(a, j, s) Q_h^k(s); a^2}{h_h^{k+d^k}(a, j, s)}
 \end{aligned}$$

where the second inequality is since  $B_h^k(s); a \leq H$ . To bound the last term, we start with a concentration bound. This allows us to substitute the indicator in Eq. (3) by its expectation (which is  $h_h^k(s); a$ ) and cancel out the denominator once. The resulting bound is:

$$\sum_{k,h;s;a} \mathbb{E}_s \sum_{q_h} \frac{h_h^{k+d^k}(a, j, s)}{h_h^k(a, j, s) + \epsilon} r_h^k(s); a \tag{6}$$

Now, the first issue we need to address is the mismatch between  $h_h^k(s)$  in the nominator and  $h_h^k(s)$  in the denominator. A similar challenge also arises in the non-delayed analysis, but in the case of delayed feedback, this requires a carefully constructed delay-adapted local bonus (defined in Eq. (4)). Note that our definition of  $h_h^k(s)$  with  $\epsilon = 3\epsilon = H$  implies that Eq. (6) is equal to  $\sum_{k,h;s} \mathbb{E}_s \sum_{q_h} h_h^k(s); i$ . Next, we

apply the value difference lemma a second time, to show that  $\text{BONUS} = \sum_{k,h;s} \mathbb{E}_s \sum_{q_h} h_h^k(s); i - \sum_{k,h;s} \mathbb{E}_s \sum_{q_h} h_h^k(s); i$ . Essentially, this means that by summing REG and BONUS, we can substitute  $h_h^k(s)$  in Eq. (6) by  $h_h^k(s)$ .

The second issue, which is unique to delayed feedback, is the mismatch between  $h_h^{k+d^k}(a, j, s)$  and  $h_h^k(a, j, s)$ . It is important to note that while in MAB  $\frac{h_h^{k+d^k}(a, j, s)}{h_h^k(a, j, s)} = \frac{h_h^k(a, j, s)}{h_h^k(a, j, s)}$  is always bounded by a constant, in MDPs this ratio can be as large as  $e^{d_{\max}}$  (see Remark B.5 in Appendix B). Thus, the standard importance-sampling estimator (with  $h_h^k(s); a$ ) will not work in this type of analysis. The main idea behind our delay-adapted estimator is that  $h_h^{k+d^k}(a, j, s) / h_h^k(a, j, s)$  which guarantees that the ratio is simply bounded by 1. Overall, we get  $\text{REG} + \text{BONUS} \leq H^2 \sum_{k,h;s;a} 1 = H^3 \text{SAK}$ , and then,  $\text{REG} + \text{BONUS} \leq H^5 K + H^3 \text{SAK}$ .

While  $r_h^k(s); a$  in our delay-adapted estimator reduces variance, it increases bias. Remarkably, this additional bias scales similarly to the DRIFT term. More specifically, for  $\text{BIAS}_1$  we first use a variant of Freedman's inequality which is highly sensitive to the estimator's variance. This brings similar issues to the ones we faced in Eq. (6), which are treated in a similar manner. Then, we show that the additional bias introduced by the ratio  $h_h^k(s); a$  scales as

$$\sum_{k,h;s;a} \mathbb{E}_s \sum_{q_h} \frac{h_h^{k+d^k}(a, j, s) (1 - r_h^k(s); a) Q_h^k(s); a}{h_h^{k+d^k}(a, j, s)} \tag{7}$$

where the inequality follows by plugging in the definition of  $r_h^k(s); a$  and some simple algebra ( $h_h^k(s); a \leq H$ ).

Utilizing the exponential weights update form, we bound the  $\ell_1$ -distance above by  $\sum_{j \in M^k} \mathbb{E}_s \sum_{q_h} \frac{h_h^{j+d^j}(a, j, s) Q_h^j(s); a}{h_h^{j+d^j}(a, j, s)}$ , where  $M^k$  is the set of episodes that their feedback arrives between episodes  $k$  and  $k + d^k$ . Then, we sum over  $k$  and apply a concentration bound over  $Q_h^k(s); a$ , which is smaller than  $Q_h^k(s); a$  in expectation (since  $h_h^k(s); a \leq 1$ ). Thus, we get that the right-hand-side (RHS) of Eq. (7) is bounded by  $H^3 \sum_k |M^k|$ , which in turn we bound by  $H^3(K + D)$  using standard delayed feedback analysis (Lemma E.10).

We can also bound the DRIFT term by the RHS of Eq. (7), up to a factor of  $H^2$  since  $\sum_j h_h^j(s); a + \sum_j h_h^j(s); a \leq H^2$ . Thus, we get that DRIFT  $\leq H^5(K + D)$  in total. The previous state-of-the-art (Jin et al., 2022) was only able to bound the DRIFT term with an additional  $\text{SA}$  factor, in part due to their complex update rule that requires solving a global optimization problem. This is a great demonstration of how a simple update rule is not only beneficial on the practical side, but also for enhanced provable guarantees. Finally,  $\text{BIAS}_2 \leq H^5$  by standard arguments for optimistic

Algorithm 2 DAPO for LinearQ-function

Initialization: Define  $\hat{Q}_h^1(a|s) = 1$  for every  $(s; a; h)$ .  
 for  $k = 1; 2; \dots; K$  do  
   Play episode  $k$  with policy  $\pi^k$ .  
   # Policy Evaluation  
   for  $j$  such that  $j + d^j = k$  do  
     Observe bandit feedback  $(r_h^j(s_h^j; a_h^j), g_{h=1}^H)$ .  
     Compute the estimated inverse covariance matrix  $\hat{\Sigma}_h^{j;+}$  via Matrix Geometric Resampling, and the estimated Q-function weights  $\hat{Q}_h^j$  defined in Eq. (8).  
     Define the delay-adapted estimator  $\hat{Q}_h^j(s; a)$  using Eq. (9), and estimate the local bonuses  $\hat{B}_h^j(s; a)$  using Algorithm 6 with respect to the local bonuses defined in Eq. (10).  
   end for  
   # Policy Update  
   Define the policy  $\pi^{k+1}$  for every  $(s; a; h)$  by:

$$\pi_h^{k+1}(a|s) = \frac{e^{\sum_{j:j+d^j=k} \hat{Q}_h^j(s;a) - \hat{B}_h^j(s;a)}}}{\sum_{a'} e^{\sum_{j:j+d^j=k} \hat{Q}_h^j(s;a) - \hat{B}_h^j(s;a)}}$$

end for

estimators. To finish the proof, sum the regret from all terms and set  $\beta = 1 - \frac{1}{SAK} + H^2(K + D)$  and  $\beta = \beta H$ .  $\square$

4. DAPO for Linear-Q

In this section we extend DAPO to linear function approximation under the LinearQ assumption (see Assumption 2.1), which generalizes Linear MDPs (Jin et al., 2020b) and in particular is much more general than the tabular setting. This enables our algorithm to scale to MDPs with a huge (possibly infinite) number of states, and gives the first regret bound for delayed feedback in non-tabular MDPs. DAPO for LinearQ (presented in Algorithm 2) follows the same framework as in Section 3. That is, in each episode there is a policy evaluation step and then a policy improvement step. Since the improvement takes the same exponential weights form, we focus on the evaluation step. Specifically, we describe the new estimator  $\hat{Q}_h^k(s; a)$  and bonus  $\hat{B}_h^k(s; a)$ , as these are the only changes compared to the tabular setting. Here we only provide sketches for the algorithm and analysis, but the full details are found in Appendix D.

Just like in the tabular case, our function estimator will simply take the original estimator from the non-delayed setting (Luo et al., 2021) and multiply it by the delay-adapted ratio. Here, the difference between our delay adaptation approach and that of Jin et al. (2022) becomes evident. While their approach of directly changing the importance-sampling weights simply does not apply anymore, our delay-adapted

ratio is easily computed locally in the current state. For completeness, we now briefly describe the estimator.

Recall that, by Assumption 2.1, the Q-function of policy  $\pi^k$  is parameterized by vectors  $\{Q_h^k\}_{h=1}^H$ , so instead of constructing an estimator  $\hat{Q}_h^k(s; a)$  for each state (which is not feasible anymore), we directly estimate  $Q_h^k$ . To that end, we first construct an estimator  $\hat{\Sigma}_h^{k;+}$  of the inverse covariance matrix  $(\Sigma_h^k + I)^{-1}$ , where  $\Sigma_h^k = E_{s;a} [g_h^k(s; a) g_h^k(s; a)^\top]$  and  $\epsilon > 0$  is an exploration parameter.  $\hat{\Sigma}_h^{k;+}$  is computed via the Matrix Geometric Resampling procedure (Neu & Olkhovskaya, 2021) which samples trajectories of the policy using the simulator. Now, the estimator  $\hat{Q}_h^k$  is defined by

$$\hat{Q}_h^k = \hat{\Sigma}_h^{k;+} \Sigma_h^k L_h^k; \tag{8}$$

where  $L_h^k$  is the cost-to-go from  $(s_h^k; a_h^k)$ . Now, to get  $\hat{Q}_h^k(s; a)$  we compute the delay adapted ratio  $\hat{r}_h^k(s; a)$  (as in Eq. (2)) and multiply it by  $\Sigma_h^k(s; a)^\top \hat{\Sigma}_h^k$ , i.e.,

$$\hat{Q}_h^k(s; a) = r_h^k(s; a) \Sigma_h^k(s; a)^\top \hat{\Sigma}_h^k. \tag{9}$$

Intuitively, this estimator is a direct generalization of the tabular importance-sampling estimator (Eq. (3)) since it corresponds to  $\frac{1}{q_h^k(s;a)+}$ , making  $\hat{Q}_h^k(s; a)$  an unbiased estimate of  $Q_h^k(s; a)$  up to  $\epsilon$  and approximation errors.

Next, we design the local bonuses  $\hat{B}_h^k(s; a)$  to go with our delay-adapted estimator. It is defined as the sum of the following local bonuses (where  $\beta_i, \gamma_i$  are parameters):

$$\begin{aligned} b_h^{k,v}(s) &= \frac{1}{v} m^{k+d^k} \sum_a r_h^k(s; a) \frac{k^{k+d^k}(a|s) k_h(s; a) k_{\Lambda_h^k}^2}{k_h(s; a) k_{\Lambda_h^k}^2} \\ b_h^{k,1}(s) &= \frac{1}{1} \sum_a r_h^k(s; a) \frac{k^{k+d^k}(a|s) k_h(s; a) k_{\Lambda_h^k}^2}{k_h(s; a) k_{\Lambda_h^k}^2} \\ b_h^{k,2}(s; a) &= \frac{1}{2} r_h^k(s; a) k_h(s; a) k_{\Lambda_h^k}^2 \\ b_h^{k,r}(s; a) &= \frac{1}{r} (1 - r_h^k(s; a)) \\ b_h^{k,f}(s) &= \frac{1}{f} \sum_a r_h^k(s; a) \frac{k^{k+d^k}(a|s) k_h(s; a) k_{\Lambda_h^k}^2}{k_h(s; a) k_{\Lambda_h^k}^2} \\ b_h^{k,g}(s; a) &= \frac{1}{g} r_h^k(s; a) k_h(s; a) k_{\Lambda_h^k}^2; \end{aligned} \tag{10}$$

where  $k_{\Lambda} = \frac{1}{x^\top A x}$  for  $x \in \mathbb{R}^n$  and  $A \in \mathbb{R}^{n \times n}$ . Each of these terms plays a different important role in the analysis.  $b_h^{k,v}(s)$  helps us control the variance of the estimator. It is inspired by Luo et al. (2021) and adapted to delay via the ratio  $r_h^k(s; a)$ .  $b_h^{k,1}(s)$  and  $b_h^{k,2}(s; a)$  help us to control the bias of the estimator. These, on the other hand, are constructed in a different manner than Luo et al. (2021). Importantly, the corresponding bonus terms in (Luo et al., 2021) might be of order  $K^{1-3}$ , while with our construction the local bonus is bounded by  $O(H^{-1} \bar{n})$ . This novel construction is what allows us to avoid dilated Bellman equations in the definition of the global bonus  $\hat{B}_h^k(s; a)$ , and by that to greatly

simplify both the algorithm and the analysis. Specifically designed to enhance exploration under delayed feedback, and in particular to control the additional bias due to the delay-adapted ratio. Finally, we add the novel terms and to ensure regret with high probability (see events and in Lemma D.3 in Appendix D). The exact role and interpretation of each of the bonus terms is further described through the main steps in the proof sketch of Theorem 4.1.

Given the local bonuses, we define to be the Q-function with respect to . However, due to the possibly infinite number of states is infeasible to compute without additional structure. Instead we compute an unbiased estimate using the simulator (see further details in Appendix D). Importantly, this estimate satisfies the Bellman equations in expectation and thus inherits some of the desired properties of Q-functions such as the validity of the value difference lemma. We note that while are defined for all states, it is sufficient to calculate them on-the-fly only over the visited states. Finally, Algorithm 6 (which we borrow from Luo et al. (2021)) that computes is not sample efficient. However, with some additional structure we can replace it with an efficient procedure recently presented by Sherman et al. (2023) and obtain the same regret (see Remark D.2).

**Theorem 4.1.** Running DAPO in a Linear adversarial MDP with  $n=K$ ,  $\epsilon = \min\{\frac{1}{10Hd_{\max}}, \frac{1}{H(K+D)^{3=4}}\}$  and access to a simulator guarantees, with probability  $1-\delta$ , that

$$R_K = \mathcal{O}(H^3 n^{5=4} K^{3=4} + H^2 D^{3=4});$$

This is the first sub-linear regret for non-tabular MDPs with delayed feedback. Moreover, like the optimal bound for tabular MDP, the delay term does not depend on the dimension  $n$ . Importantly, our analysis is relatively simple even compared to the non-delayed case. By that we lay solid foundations for improved regret bounds in future work, and manage to bound the regret with high probability and not just in expectation. One significant difference between the tabular case and Linear is that now the estimator might be negative (specifically,  $H=$ ). This is not a problem in the non-delayed setting which does not have a drift term, and in fact, with proper hyper-parameter tuning we get the same  $\mathcal{O}(H^2 n^{2=3} K^{2=3})$  bound of Luo et al. (2021) without delays. However, in the presence of delays this issue induces new challenges and requires a more involved analysis and algorithmic design (and leads to worse regret).

**Proof sketch of Theorem 4.** We start by decomposing the regret as in Eq. (5). To bound  $\text{REG}$  we can no longer apply Corollary E.7 like we did in the tabular case, because it heavily relies on the losses being not too negative (now they might be  $\mathcal{O}(H=)$ ). Instead, we prove a novel bound

(Lemma E.8) that bounds  $\text{REG}$ , for sufficiently small  $\epsilon$ , by  $H^2 \sum_{k,h} E_{s, q_h; a} m^{k+d^k} (\hat{Q}_h^k(s; a) - \hat{B}_h^k(s; a))^2$ ;

where  $m^k = \sum_j |j| : j + d^j = k$ . Next, follow similar steps to Theorem 3.1. Specifically, we use  $\hat{B}_h^k(s; a) = H^2 \bar{n}$ , apply a concentration bound  $\mathcal{O}(\hat{Q}_h^k(s; a)^2)$  around its expectation and further bound the expectation. This allows us to show that:  $\text{REG} \leq H + H^5 n K + H^2 \sum_{k,h} E_{s, q_h; a} m^{k+d^k} r_h^k(s; a) k_h(s; a) k_{\Lambda_h^{k;+}}^2$ .

Now we once again face the issue that the expectation is taken over states generated by  $\hat{Q}$  and actions generated by  $k+d^k$ , while  $\Lambda_h^{k;+}$  is constructed from trajectories generated by  $k$ . Remarkably, our technique from the tabular case extends naturally to Linear. Since  $\hat{B}_h^k(s; a)$  satisfies the Bellman equations in expectation, we can use the value difference lemma to show that  $\text{BONUS} = \sum_{k,h} E_{s; a, q_h} [b_h^k(s; a)] - \sum_{k,h} E_{s; a, q_h} [b_h^k(s; a)]$ .

Recall that  $\hat{Q}_h^k(s; a)$  is the sum of the  $6$  local bonuses defined in Eq. (10), so we can write  $\text{BONUS} = \text{BONUS}_v + \text{BONUS}_1 + \text{BONUS}_2 + \text{BONUS}_3 + \text{BONUS}_4 + \text{BONUS}_5 + \text{BONUS}_6$ , where each term corresponds to its local bonus. We set  $v = H^2$  to get that  $(\cdot) = \sum_{k,h} E_{s, q_h} [b_h^{k;v}(s)]$ , so naturally  $(\cdot) + \text{BONUS}_v$  is bounded by,

$$H^2 \sum_{k,h} E_{s, q_h^k; a} m^{k+d^k} r_h^k(s; a) k_h(s; a) k_{\Lambda_h^{k;+}}^2 + H^2 \sum_{k,h} E_{s, q_h^k; a} m^{k+d^k} k_h(s; a) k_{\Lambda_h^{k;+}}^2; \quad (11)$$

where the last step is due to our delay-adapted ratio, demonstrating its power compared to directly changing the importance-sampling weights (Jin et al., 2022). It lets us adjust the expectation to be over trajectories sampled with  $k$ , so it is aligned with the construction  $\Delta_h^{k;+}$  via Matrix Geometric Resampling. To further bound Eq. (11), we may now utilize standard techniques from non-delayed analysis (e.g., Jin et al. (2020b); Luo et al. (2021)). We plug in the definition of the matrix norm, then we can consider its trace and use its linearity and invariance under cyclic permutations. This enables us to bound the expectation in Eq. (11) by  $\text{tr} \Lambda_h^{k;+} E_{s; a, q_h^k} \langle \hat{Q}_h^k(s; a) \hat{Q}_h^k(s; a) \rangle = \text{tr} \Lambda_h^{k;+} \frac{1}{k} \cdot$ . Finally, since  $\Lambda_h^{k;+}$  approximates  $(\frac{k}{h} + I)^{-1}$ , the last term is approximately bounded by  $\frac{1}{k}$ . Thus, we get that  $\text{BONUS}_v + \text{REG} \leq H + H^5 n K$  because  $\sum_k m^{k+d^k} \leq K$ .

For the analysis of  $\text{BIAS}_1$ , we first show it is mainly bounded by two terms: the first comes from the standard estimator while the second is the additional bias due to the delay-

adapted ratio  $\hat{r}_h^k(s; a)$ . That is, we bound  $\text{BIAS}_1$  by:

$$H^P \frac{X}{n} \mathbb{E}_{s, q_h; a} r_h^k(s; a) k_h(s; a) k_h^{\wedge k; +} \quad (12)$$

$$+ H^P \frac{X}{n} \mathbb{E}_{s, q_h; a} (1 - r_h^k(s; a)) : \quad (13)$$

Once again, with the proper tuning  $\beta = H^P \frac{X}{n}$ , we can get (12) =  $\mathbb{E}_{s, q_h} [b_h^{k; 1}(s)]$ , and then combine with the corresponding BONUS term  $\text{BONUS}_1$ , while utilizing the delay-adapted ratio. This gives (12) +  $\text{BONUS}_1 = H^P \frac{X}{n} \mathbb{E}_{s; a} q_h^k k_h(s; a) k_h^{\wedge k; +} - H^2 n K$ .

For Eq. (13), we first bound  $\mathbb{E}_a (1 - r_h^k(s; a))$

$k_h^{k+d^k}(j_s) - k_h^k(j_s) k_1$  and then follow similar arguments to the tabular case regarding the multiplicative weights update form. The main difference is that now  $\hat{r}_h^k(s; a)$  can be negative, resulting in weaker guarantees (see more details in Appendix D.6). Overall, we get (13) =  $-H^3 \frac{P}{n} (K + D)$ .

The analysis of  $\text{BIAS}_2$  is very different from both the tabular case and the non-delayed Linear (Luo et al., 2021). This is mainly for two reasons: First, in the tabular case, the added bias makes the estimator  $\hat{r}_h^k(s; a)$  optimistic, but this is no longer the case in Linear since now the estimator might also be negative. Second,  $\text{BIAS}_2$  contains the inner product with  $k_h(j_s)$  which cannot be aligned with the denominator of the delay-adapted ratio  $\hat{r}_h^k(s; a)$ . Thus, we need a novel more involved analysis for  $\text{BIAS}_2$ .

We start in a similar way to  $\text{BIAS}_1$ , and bound  $\text{BIAS}_2$  by:

$$H^P \frac{X}{n} \mathbb{E}_{s; a} q_h r_h^k(s; a) k_h(s; a) k_h^{\wedge k; +} \quad (14)$$

$$+ H^P \frac{X}{n} \mathbb{E}_{s; a} q_h (1 - r_h^k(s; a)) : \quad (15)$$

We handle Eq. (14) like Eq. (12), so for  $\beta = H^P \frac{X}{n}$  we get that (14) +  $\text{BONUS}_2 = H^2 n K$ . Term (15) on the other hand might be of order  $\beta^2$  due to mismatch between  $k_h(a_j s)$  in the expectation and  $\max_h^{k+d^k}(a_j s); k_h(a_j s) g$  in the denominator of  $r_h^k(s; a)$ . To address that, we design the novel bonus term  $k_h^{k; r}(s; a)$ . Summing Eq. (15) with  $\text{BONUS}_2$  essentially allows us to substitute  $k_h(a_j s)$  by  $k_h^k(a_j s)$ , which gives: (14) +  $\text{BONUS}_2 = H^P \frac{X}{n} \mathbb{E}_{s, q_h} k_h^{k+d^k}(j_s) - k_h^k(j_s) k_1 - H^3 \frac{P}{n} (K + D)$ .

Finally, DRIFT is also bounded by the  $\beta$ -distance above and further by  $-H^4 (K + D)$ . Summing the regret from all terms and optimizing over  $\beta$  completes the proof.  $\square$

## 5. Delay-Adapted PPO and Experiments

In this section we show how our generic delay adaptation method extends to state-of-the-art deep RL methods, and demonstrate its great potential through simple experiments on popular MuJoCo environments (Todorov et al., 2012). We note that here we follow the standard convention that use rewards in deep RL rather than costs as in the rest of the paper. Due to lack of space, here we only provide an overview of the method and main experiment. For additional experiments and full implementation details see Appendix F.

The highly successful TRPO algorithm (Schulman et al., 2015) is a deep RL policy-gradient method that builds on the same PO principles discussed in this paper. Specifically, it follows the Policy Iteration paradigm with a ‘‘soft’’ policy improvement step, where the policy is now approximated by a Deep Neural Network with parameters  $\theta$ . While our update step (Eq. (1)) is equivalent to maximizing an objective with a KL-regularization term (Shani et al., 2020a), TRPO replaces it by a constraint which results in maximizing the objective  $L_{\text{TRPO}}^k(\theta) = \sum_{h=1}^H \frac{(a_h j_s \theta)}{k(a_h j_s \theta)} \hat{A}_h$  subject to a constraint that keeps the new and the old policies close in terms of KL-divergence. Here  $\hat{A}_h$  is an estimate of the advantage function which replaces  $Q$  function in our formulation to further reduce variance (Sutton et al., 1999).

While successful, TRPO’s constrained optimization is computationally expensive. Thus, PPO (Schulman et al., 2017) removes the explicit constraint and replaces it with a sophisticated clipping technique that allows to keep strong empirical performance while only optimizing the following (non-constrained) objective:

$$L^k(\theta) = \sum_{h=1}^H \min \left( g_h^k(\theta) \hat{A}_h; \text{clip}_1 \left( g_h^k(\theta) \hat{A}_h^o \right) \right); \quad (16)$$

where  $g_h^k(\theta) = \frac{(a_h j_s \theta)}{k(a_h j_s \theta)}$  and  $\text{clip}_1(x)$  clips  $x$  between  $1 - \beta$  and  $1 + \beta$ . This objective essentially zeros the gradient whenever the policy changes too much, and thus replaces the need for an explicit constraint (see more details in (Schulman et al., 2017)). Next, we adapt PPO to delayed feedback.

With delayed feedback, the trajectory that arrives at time  $t$  was generated using policy  $\pi^{k^d}$ . Thus, a naive adaptation would be to optimize Eq. (16) but replacing  $g^k$  with  $g^{k^d}$ . We will refer to this algorithm as Delayed PPO (DPPO). As this paper shows, DPPO is likely to suffer from large variance which can be reduced when multiplying the objective by the delay-adapted ratio  $\hat{r}_h^k(s; a)$ . The result is our novel Delay-Adapted PPO (DAPPO) which optimizes:

$$L_{\text{DA}}^k(\theta) = \sum_{h=1}^H \min \left( R_h^k(\theta) \hat{A}_h; \text{clip}_1 \left( R_h^k(\theta) \hat{A}_h^o \right) \right);$$

for  $R_h^k(\theta) = \frac{(a_h j_s \theta)}{\max \left( k^d(a_h j_s \theta); k(a_h j_s \theta) g \right)}$ . Another alter-

native, which we call Non-Delayed PPO (NDPPO), is to use the original PPO and ignore the fact that feedback is



Figure 1. Training curves: DAPPO vs DPPO. Plots show average reward and std over 5 seeds. x-axis is number of timesteps up to 5M.

Figure 2. Training curves with different xed delay length: DAPPO vs DPPO with different delay, alongside PPO without delays. Plots show average reward and std over 5 seeds. x-axis is number of timesteps up to 5M.

delayed. This results in a highly unstable algorithm which is likely to suffer from large bias due to the mismatch between the policy  $\pi^k$  that generated the trajectory and the policy  $\pi^d$  which is used to re-weight the estimator.

Fig. 1 compares the performance of DPPO and DAPPO over 8 MuJoCo environments. We use a xed delay of 10<sup>6</sup> timesteps while the total number of timesteps is 10<sup>7</sup>. Results are averaged over 5 runs and the shaded areas around the curves indicate standard deviation. Note that the only difference between the algorithms is the objective, we did not tune hyper-parameters or modify network architecture. DAPPO outperforms DPPO in at least 4 environments and is on par with DPPO in the rest. The only exception is InvertedDoublePendulum, but it is important to note that this environment is extremely noisy.

Fig. 2 compares the training curves of DPPO vs DAPPO in the SWIMMER environment (for more environments see Appendix F.3) with different delays in {10000, 25000, 50000, 75000, 100000}, alongside the training curve of PPO without delay. As expected, when the delay is relatively small (e.g. 10000), there is no significant difference between learning with or without delayed feedback. As the delay becomes larger, the performance of all algorithms drops (but at different rates).

These empirical results support our claim that handling

delays via the delay-adapted ratio extends naturally beyond the tabular and Linear Q settings to practical deep function approximation. Surprisingly, even in this simple case, delays cause significant drop in performance which demonstrates the great importance of delay-adapted algorithms. Our novel method makes a significant step towards practical deep RL algorithms that are robust to delayed feedback. Finally, we note that NDPPO is omitted from the graphs because it does not converge (as expected). Instead, it oscillates between high and low reward (see Appendix F for more details).

## 6. Future Work

We leave a few open questions for future work. In the tabular case, it still remains unclear what is the optimal dependency under delayed feedback in terms of the horizon. Another future direction is to further improve our results in the Linear Function Approximation setting and extend them to the case where a simulator is unavailable. This is in particular important in light of very recent advancement in the non-delayed setting (Sherman et al., 2023; Dai et al., 2023) which significantly improve Luo et al. (2021). Finally, our experiment demonstrate the potential that delay-adaptation methods have for deep RL applications. However, a much more thorough empirical study needs to be done in order to fully understand the implications of these methods on deep RL.

## References

- Abbasi-Yadkori, Y., Bartlett, P., Bhatia, K., Lazic, N., Szepesvari, C., and Weisz, G. Politex: Regret bounds for policy iteration using expert prediction. *International Conference on Machine Learning*, pp. 3692–3702. PMLR, 2019.
- Agarwal, A. and Duchi, J. C. Distributed delayed stochastic optimization. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pp. 5451–5452. IEEE, 2012.
- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. Optimality and approximation with policy gradient methods in markov decision processes. *Conference on Learning Theory*, pp. 64–66. PMLR, 2020.
- Ayoub, A., Jia, Z., Szepesvari, C., Wang, M., and Yang, L. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pp. 463–474. PMLR, 2020.
- Azar, M. G., Osband, I., and Munos, R. Minimax regret bounds for reinforcement learning. *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pp. 263–272. JMLR. org, 2017.
- Beygelzimer, A., Langford, J., Li, L., Reyzin, L., and Schapire, R. Contextual bandit algorithms with supervised learning guarantees. *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 19–26. JMLR Workshop and Conference Proceedings, 2011.
- Bistriz, I., Zhou, Z., Chen, X., Bambos, N., and Blanchet, J. Online  $\epsilon_3$  learning in adversarial bandits with delayed feedback. In *Advances in Neural Information Processing Systems*, pp. 11349–11358, 2019.
- Bistriz, I., Zhou, Z., Chen, X., Bambos, N., and Blanchet, J. No discounted-regret learning in adversarial bandits with delays. *arXiv preprint arXiv:2103.04550*, 2021.
- Cai, Q., Yang, Z., Jin, C., and Wang, Z. Provably efficient exploration in policy optimization. *International Conference on Machine Learning*, pp. 1283–1294. PMLR, 2020.
- Cesa-Bianchi, N., Gentile, C., Mansour, Y., and Minora, A. Delay and cooperation in nonstochastic bandits. In *Conference on Learning Theory*, pp. 605–622, 2016.
- Cesa-Bianchi, N., Gentile, C., and Mansour, Y. Nonstochastic bandits with composite anonymous feedback. *Conference On Learning Theory*, pp. 750–773, 2018.
- Cesa-Bianchi, N., Gentile, C., and Mansour, Y. Delay and cooperation in nonstochastic bandits. *The Journal of Machine Learning Research*, 20(1):613–650, 2019.
- Changuel, N., Sayadi, B., and Kieffer, M. Online learning for qoe-based video streaming to mobile receivers. In *2012 IEEE Globecom Workshops*, pp. 1319–1324. IEEE, 2012.
- Chen, B., Xu, M., Liu, Z., Li, L., and Zhao, D. Delay-aware multi-agent reinforcement learning. *arXiv preprint arXiv:2005.05441*, 2020a.
- Chen, L. and Luo, H. Finding the stochastic shortest path with low regret: The adversarial cost and unknown transition case. *arXiv preprint arXiv:2102.05284*, 2021.
- Chen, L., Luo, H., and Wei, C.-Y. Minimax regret for stochastic shortest path with adversarial costs and known transition. *arXiv preprint arXiv:2012.04053*, 2020b.
- Chen, L., Jafarnia-Jahromi, M., Jain, R., and Luo, H. Implicit finite-horizon approximation and efficient optimal algorithms for stochastic shortest path. *Advances in Neural Information Processing Systems*, 2021.
- Chen, L., Jain, R., and Luo, H. Improved no-regret algorithms for stochastic shortest path with linear mdp. In *International Conference on Machine Learning*, pp. 3204–3245. PMLR, 2022a.
- Chen, L., Luo, H., and Rosenberg, A. Policy optimization for stochastic shortest path. In Loh, P. and Raginsky, M. (eds.), *Conference on Learning Theory*, 2-5 July 2022, London, UK, volume 178 of *Proceedings of Machine Learning Research*, pp. 982–1046. PMLR, 2022b.
- Cohen, A., Daniely, A., Drori, Y., Koren, T., and Schain, M. Asynchronous stochastic optimization robust to arbitrary delays. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 9024–9035, 2021a.
- Cohen, A., Efroni, Y., Mansour, Y., and Rosenberg, A. Minimax regret for stochastic shortest path. *Advances in Neural Information Processing Systems*, 34, 2021b.
- Dai, Y., Luo, H., and Chen, L. Follow-the-perturbed-leader for adversarial markov decision processes with bandit feedback. *arXiv preprint arXiv:2205.13451*, 2022.
- Dai, Y., Luo, H., Wei, C.-Y., and Zimmert, J. Reduced regret for adversarial mdps with linear function approximation. *arXiv preprint arXiv:2301.12942*, 2023.

- Dann, C., Lattimore, T., and Brunskill, E. Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 5713–5723, 2017.
- Derman, E., Dalal, G., and Mannor, S. Acting in delayed environments with non-stationary markov policies. *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.
- Dudik, M., Hsu, D., Kale, S., Karampatziakis, N., Langford, J., Reyzin, L., and Zhang, T. Efficient optimal learning for contextual bandits. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, pp. 169–178, 2011.
- Efroni, Y., Merlis, N., Ghavamzadeh, M., and Mannor, S. Tight regret bounds for model-based reinforcement learning with greedy policies. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada*, pp. 12203–12213, 2019.
- Even-Dar, E., Kakade, S. M., and Mansour, Y. Online markov decision processes. *Mathematics of Operations Research* 34(3):726–736, 2009.
- 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.
- Gael, M. A., Vernade, C., Carpentier, A., and Valko, M. Stochastic bandits with arm-dependent delays. *International Conference on Machine Learning*, pp. 3348–3356. PMLR, 2020.
- Gu, S., Holly, E., Lillicrap, T., and Levine, S. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. *2017 IEEE international conference on robotics and automation (ICRA)*, pp. 3389–3396. IEEE, 2017.
- Gyorgy, A. and Joulani, P. Adapting to delays and data in adversarial multi-armed bandits. *International Conference on Machine Learning*, pp. 3988–3997. PMLR, 2021.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pp. 1861–1870, 2018.
- He, J., Zhou, D., and Gu, Q. Near-optimal policy optimization algorithms for learning adversarial linear mixture mdps. In Camps-Valls, G., Ruiz, F. J. R., and Valera, I. (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 4259–4280. PMLR, 28–30 Mar 2022.
- Howson, B., Pike-Burke, C., and Filippi, S. Delayed feedback in episodic reinforcement learning. *arXiv preprint arXiv:2111.07615*, 2021.
- Howson, B., Pike-Burke, C., and Filippi, S. Delayed feedback in generalised linear bandits revisited. *arXiv preprint arXiv:2207.10786*, 2022.
- Ito, S., Hatano, D., Sumita, H., Takemura, K., Fukunaga, T., Kakimura, N., and Kawarabayashi, K.-I. Delay and cooperation in nonstochastic linear bandits. *Advances in Neural Information Processing Systems* 33:4872–4883, 2020.
- Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research* 11(4), 2010.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. Is q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pp. 4863–4873, 2018.
- Jin, C., Jin, T., Luo, H., Sra, S., and Yu, T. Learning adversarial markov decision processes with bandit feedback and unknown transition. In *International Conference on Machine Learning*, pp. 4860–4869. PMLR, 2020a.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pp. 2137–2143, 2020b.
- Jin, T. and Luo, H. Simultaneously learning stochastic and adversarial episodic mdps with known transition. *Advances in neural information processing systems* 2020, 2020.
- Jin, T., Huang, L., and Luo, H. The best of both worlds: stochastic and adversarial episodic mdps with unknown transition. *Advances in Neural Information Processing Systems* 2021.
- Jin, T., Lancewicki, T., Luo, H., Mansour, Y., and Rosenberg, A. Near-optimal regret for adversarial mdp with delayed bandit feedback. *arXiv preprint arXiv:2201.13172*, 2022.
- Kakade, S. and Langford, J. Approximately optimal approximate reinforcement learning. In *Proc. 19th International Conference on Machine Learning*, pp. 405–412, 2002.

- Kakade, S. M. A natural policy gradient. *Advances in neural information processing systems* 14:1531–1538, 2001.
- Katsikopoulos, K. V. and Engelbrecht, S. E. Markov decision processes with delays and asynchronous cost collection. *IEEE transactions on automatic control* 48(4): 568–574, 2003.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lancewicki, T., Segal, S., Koren, T., and Mansour, Y. Stochastic multi-armed bandits with unrestricted delay distributions. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, pp. 5969–5978. PMLR, 2021.
- Lancewicki, T., Rosenberg, A., and Mansour, Y. Cooperative online learning in stochastic and adversarial mdps. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvári, C., Niu, G., and Sabato, S. (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 11918–11968. PMLR, 2022a.
- Lancewicki, T., Rosenberg, A., and Mansour, Y. Learning adversarial markov decision processes with delayed feedback. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7281–7289, 2022b.
- Levine, S. and Koltun, V. Guided policy search. *International conference on machine learning*, pp. 1–9. PMLR, 2013.
- Levine, S., Finn, C., Darrell, T., and Abbeel, P. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research* 17(1):1334–1373, 2016.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Liu, S., Wang, X., and Liu, P. X. Impact of communication delays on secondary frequency control in an islanded microgrid. *IEEE Transactions on Industrial Electronics* 62(4):2021–2031, 2014.
- Luo, H., Wei, C.-Y., and Lee, C.-W. Policy optimization in adversarial mdps: Improved exploration via diluted bonuses. *Advances in Neural Information Processing Systems* 34, 2021.
- Mahmood, A. R., Korenkevych, D., Komer, B. J., and Bergstra, J. Setting up a reinforcement learning task with a real-world robot. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4635–4640. IEEE, 2018.
- Masoudian, S., Zimmert, J., and Seldin, Y. A best-of-both-worlds algorithm for bandits with delayed feedback. *arXiv preprint arXiv:2206.14906*, 2022.
- Min, Y., He, J., Wang, T., and Gu, Q. Learning stochastic shortest path with linear function approximation. In *International Conference on Machine Learning*, pp. 15584–15629. PMLR, 2022.
- Neu, G. Explore no more: Improved high-probability regret bounds for non-stochastic bandits. *Advances in Neural Information Processing Systems* 28:3168–3176, 2015.
- Neu, G. and Olkhovskaya, J. Online learning in mdps with linear function approximation and bandit feedback. *Advances in Neural Information Processing Systems* 34:10407–10417, 2021.
- Neu, G., György, A., and Szepesvári, C. The online loop-free stochastic shortest-path problem. *COLT 2010 - The 23rd Conference on Learning Theory, Haifa, Israel, June 27-29, 2010*, pp. 231–243, 2010a.
- Neu, G., György, A., and Szepesvári, C. The online loop-free stochastic shortest-path problem. *Conference on Learning Theory (COLT)*, pp. 231–243, 2010b.
- Neu, G., György, A., and Szepesvári, C. The adversarial stochastic shortest path problem with unknown transition probabilities. In *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 805–813, 2012.
- Neu, G., György, A., Szepesvári, C., and Antos, A. Online Markov Decision Processes under bandit feedback. *IEEE Trans. Automat. Contr* 59(3):676–691, 2014.
- Pike-Burke, C., Agrawal, S., Szepesvári, C., and Grunewalder, S. Bandits with delayed, aggregated anonymous feedback. *International Conference on Machine Learning*, pp. 4105–4113. PMLR, 2018.
- Quanrud, K. and Khashabi, D. Online learning with adversarial delays. *Advances in neural information processing systems* 28:1270–1278, 2015.
- Rafan, A., Hill, A., Gleave, A., Kanervisto, A., Ernestus, M., and Dormann, N. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research* 2021.
- Rosenberg, A. and Mansour, Y. Online stochastic shortest path with bandit feedback and unknown transition function. In *Advances in Neural Information Processing Systems*, pp. 2209–2218, 2019a.
- Rosenberg, A. and Mansour, Y. Online convex optimization in adversarial markov decision processes. In *International Conference on Machine Learning*, pp. 5478–5486. PMLR, 2019b.

- Rosenberg, A. and Mansour, Y. Oracle-efficient regret minimization in factored mdps with unknown structure. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W. (eds) *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 11148–11159, 2021a.
- Rosenberg, A. and Mansour, Y. Stochastic shortest path with adversarially changing costs. In Zhou, Z. (ed.) *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pp. 2936–2942. ijcai.org, 2021b.
- Rosenberg, A., Cohen, A., Mansour, Y., and Kaplan, H. Near-optimal regret bounds for stochastic shortest path. In *International Conference on Machine Learning*, pp. 8210–8219. PMLR, 2020.
- Schuitema, E., Boniu, L., Babuska, R., and Jonker, P. Control delay in reinforcement learning for real-time dynamic systems: a memoryless approach. *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 3226–3231. IEEE, 2010.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization. *International conference on machine learning*, pp. 1889–1897, 2015.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Shani, L., Efroni, Y., and Mannor, S. Adaptive trust region policy optimization: Global convergence and faster rates for regularized mdps. *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pp. 5668–5675. AAAI Press, 2020a.
- Shani, L., Efroni, Y., Rosenberg, A., and Mannor, S. Optimistic policy optimization with bandit feedback. In *International Conference on Machine Learning*, pp. 8604–8613. PMLR, 2020b.
- Sherman, U., Koren, T., and Mansour, Y. Improved regret for efficient online reinforcement learning with linear function approximation. *arXiv preprint arXiv:2301.13087*, 2023.
- Sutton, R. S. and Barto, A. *Reinforcement learning: An introduction*. MIT press, 2018.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 2, 1999.
- Tarbouriech, J., Garcelon, E., Valko, M., Pirotta, M., and Lazaric, A. No-regret exploration in goal-oriented reinforcement learning. *International Conference on Machine Learning*, pp. 9428–9437. PMLR, 2020.
- Tarbouriech, J., Zhou, R., Du, S. S., Pirotta, M., Valko, M., and Lazaric, A. Stochastic shortest path: Minimax, parameter-free and towards horizon-free regret. *Advances in Neural Information Processing Systems*, 34, 2021.
- Thune, T. S., Cesa-Bianchi, N., and Seldin, Y. Nonstochastic multiarmed bandits with unrestricted delays. *Advances in Neural Information Processing Systems*, pp. 6541–6550, 2019.
- Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033. IEEE, 2012.
- Tomar, M., Shani, L., Efroni, Y., and Ghavamzadeh, M. Mirror descent policy optimization. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*, 2022.
- Van Der Hoeven, D. and Cesa-Bianchi, N. Nonstochastic bandits and experts with arm-dependent delays. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022.
- Vernade, C., Cappé, O., and Perchet, V. Stochastic bandit models for delayed conversions. *Conference on Uncertainty in Artificial Intelligence*, 2017.
- Vernade, C., Carpentier, A., Lattimore, T., Zappella, G., Ermis, B., and Brueckner, M. Linear bandits with stochastic delayed feedback. *International Conference on Machine Learning*, pp. 9712–9721. PMLR, 2020.
- Vial, D., Parulekar, A., Shakkottai, S., and Srikant, R. Regret bounds for stochastic shortest path problems with linear function approximation. *International Conference on Machine Learning*, pp. 22203–22233. PMLR, 2022.
- Walsh, T. J., Nouri, A., Li, L., and Littman, M. L. Learning and planning in environments with delayed feedback. *Autonomous Agents and Multi-Agent Systems*, 13(1):83, 2009.
- Wei, C.-Y., Jahromi, M. J., Luo, H., and Jain, R. Learning in finite-horizon average-reward mdps with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pp. 3007–3015. PMLR, 2021.

- Yang, L. and Wang, M. Sample-optimal parametric  $q$ -learning using linearly additive features. *International Conference on Machine Learning*, pp. 6995–7004. PMLR, 2019.
- Zanette, A., Brandfonbrener, D., Brunskill, E., Pirotta, M., and Lazaric, A. Frequentist regret bounds for randomized least-squares value iteration. *International Conference on Artificial Intelligence and Statistics*, pp. 1954–1964, 2020a.
- Zanette, A., Lazaric, A., Kochenderfer, M., and Brunskill, E. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning* pp. 10978–10989. PMLR, 2020b.
- Zhou, D. and Gu, Q. Computationally efficient horizon-free reinforcement learning for linear mixture mdp. *arXiv preprint arXiv:2205.11507*, 2022.
- Zhou, D., Gu, Q., and Szepesvari, C. Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. *Conference on Learning Theory* pp. 4532–4576. PMLR, 2021.
- Zhou, Z., Xu, R., and Blanchet, J. Learning in generalized linear contextual bandits with stochastic delays. In *Advances in Neural Information Processing Systems* pp. 5197–5208, 2019.
- Zimin, A. and Neu, G. Online learning in episodic markovian decision processes by relative entropy policy search. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States* pp. 1583–1591, 2013.
- Zimmert, J. and Seldin, Y. An optimal algorithm for adversarial bandits with arbitrary delays. *International Conference on Artificial Intelligence and Statistics*, pp. 3285–3294. PMLR, 2020.

# Appendix

---

|     |                                                                                         |    |
|-----|-----------------------------------------------------------------------------------------|----|
| A   | Related Work                                                                            | 16 |
| B   | Delay-Adapted Policy Optimization for (Tabular) Adversarial MDP with Known Transition   | 18 |
| B.1 | The good event                                                                          | 18 |
| B.2 | Proof of the main theorem                                                               | 20 |
| B.3 | Bound on $\mathbb{B}AS_1$                                                               | 21 |
| B.4 | Bound on $\mathbb{B}ONUS$                                                               | 22 |
| B.5 | Bound on $\mathbb{R}EG$                                                                 | 23 |
| B.6 | Bound on $\mathbb{D}RIFT$                                                               | 23 |
| C   | Delay-Adapted Policy Optimization for (Tabular) Adversarial MDP with Unknown Transition | 26 |
| C.1 | The good event                                                                          | 27 |
| C.2 | Proof of the main theorem                                                               | 29 |
| C.3 | Bound on $\mathbb{B}AS_1$                                                               | 30 |
| C.4 | Bound on $\mathbb{B}ONUS$                                                               | 31 |
| C.5 | Bound on $\mathbb{R}EG$                                                                 | 32 |
| D   | Delay-Adapted Policy Optimization for Adversarial MDP with Linear Q-function            | 33 |
| D.1 | The good event                                                                          | 34 |
| D.2 | Proof of the main theorem                                                               | 36 |
| D.3 | Bound on $\mathbb{B}AS_1$                                                               | 37 |
| D.4 | Bound on $\mathbb{B}AS_2$                                                               | 39 |
| D.5 | Bound on $\mathbb{R}EG$                                                                 | 40 |
| D.6 | Bound on $\mathbb{D}RIFT$                                                               | 42 |
| D.7 | Bound on $\mathbb{B}ONUS$                                                               | 43 |
| E   | Auxiliary Lemmas                                                                        | 46 |
| F   | DAPPO Implementation Details and Additional Experiments                                 | 51 |
| F.1 | DAPPO Implementation Details                                                            | 51 |
| F.2 | Additional Experiments – Instability of NDPPPO                                          | 52 |
| F.3 | Additional Experiments – Drop in Performance as Delay Length Increases                  | 52 |

---

## A. Related Work

In this section we provide a full review of the literature related to regret minimization in adversarial MDP with delayed feedback. For completeness, we include topics that are not directly related to this paper.

**Delays in RL without Regret Analysis.** Delays were studied in the practical RL literature (Schuitema et al., 2010; Liu et al., 2014; Changuel et al., 2012; Mahmood et al., 2018; Derman et al., 2021), but this is not related the topic of this paper. In the theory literature, most previous work (Katsikopoulos & Engelbrecht, 2003; Walsh et al., 2009) considered delays in the observation of the state. That is, when the agent takes an action she is not certain what is the current state, and will only observe it in delay. This setting is much more related to partially observable MDPs (POMDPs) and motivated by scenarios like robotics system delays. Unfortunately, even planning is computationally hard (exponential in the delay for delayed state observability (Walsh et al., 2009)). The topic studied in this paper (i.e., delayed feedback) is inherently different, and is motivated by settings like recommendation systems. Importantly, unlike delayed state observability, it is not computationally hard to handle delayed feedback. The challenges of delayed feedback are very different than the ones of delayed state observability, and include policy updates that occur in delay and exploration without observing feedback (Lancewicki et al., 2022b).

**Delays in RL with Regret Analysis.** This line of work is related to this paper the most. Howson et al. (2021) studied delayed feedback in stochastic MDPs, and assume that the delays are also stochastic, i.e., sampled i.i.d from a fixed (unknown) distribution. This is a restrictive assumption since it does not allow dependencies between costs and delays that are very common in practice. In adversarial MDPs, delayed feedback was first studied by Lancewicki et al. (2022b). They proposed Policy Optimization algorithms that handle delays, but focused on the case of full-information feedback where the agent observes the entire cost function in the end of the episode instead of bandit feedback (where the agent observes only costs along its trajectory). Full-information feedback is not a realistic assumption in most applications, and for bandit feedback they only prove sub-optimal regret bounds of  $\tilde{O}(D)^{2-3}$  (ignoring dependencies  $B; A; H$ ). Later, Jin et al. (2022) managed to achieve a near-optimal regret bound of  $\tilde{O}(SAK + (HSA)^{1-4}H^2D)$  for the case of known transition function and  $\tilde{O}(H^2S^2AK + (HSA)^{1-4}H^2D)$  for the case of unknown transitions. However, their algorithm is based on the O-REPS method (Zimin & Neu, 2013) which requires solving a computationally expensive global optimization problem and cannot be extended to function approximation. Recently, Dai et al. (2022) showed that delayed feedback in adversarial MDPs can also be dealt with using Follow-The-Perturbed-Leader (FTPL) algorithms. The efficiency of FTPL algorithms is similar to Policy Optimization, but their regret bound is only  $\tilde{O}(H^2S^2AK + HSA^2H^2D)$ .

**Delays in multi-arm bandit (MAB).** Delays were extensively studied in MAB and online optimization both in the stochastic setting (Dudik et al., 2011; Agarwal & Duchi, 2012; Vernade et al., 2017; 2020; Pike-Burke et al., 2018; Cesa-Bianchi et al., 2018; Zhou et al., 2019; Gael et al., 2020; Lancewicki et al., 2021; Cohen et al., 2021a; Howson et al., 2022), and the adversarial setting (Quanrud & Khashabi, 2015; Cesa-Bianchi et al., 2016; Thune et al., 2019; Bistritz et al., 2019; Zimmert & Seldin, 2020; Ito et al., 2020; Gyorgy & Joulani, 2021; Van Der Hoeven & Cesa-Bianchi, 2022; Masoudian et al., 2022). However, as discussed in (Lancewicki et al., 2022b), delays introduce new challenges in MDPs that do not appear in MAB.

**Regret minimization in Tabular RL.** There exists a rich literature on regret minimization in tabular MDPs. In the stochastic case, the algorithms are mainly built on optimism in face of uncertainty approach (Jaksch et al., 2010; Azar et al., 2017; Dann et al., 2017; Jin et al., 2018; Efroni et al., 2019; Tarbouriech et al., 2020; 2021; Rosenberg & Mansour, 2021a; Rosenberg et al., 2020; Cohen et al., 2021b; Chen et al., 2021). In the adversarial case, while a few algorithms use FTPL (Neu et al., 2012; Dai et al., 2022), most algorithms are based on either the O-REPS method (Zimin & Neu, 2013; Rosenberg & Mansour, 2019b;a; 2021b; Jin et al., 2020a; 2021; Jin & Luo, 2020; Lancewicki et al., 2022a; Chen et al., 2020b; Chen & Luo, 2021) or on Policy Optimization (Even-Dar et al., 2009; Neu et al., 2010a;b; 2014; Shani et al., 2020b; Luo et al., 2021; Chen et al., 2022b). Note that regret minimization in standard episodic MDPs is a special case of the model considered in this paper where  $\delta_k = 0$  for every episode  $k$ .

**Regret minimization in RL with Linear Function Approximation.** In recent years the literature on regret minimization in RL has expanded to linear function approximation. While in the stochastic case algorithms are still based on optimism (Jin et al., 2020b; Yang & Wang, 2019; Zanette et al., 2020a;b; Ayoub et al., 2020; Zhou et al., 2021; Zhou & Gu, 2022; Vial et al., 2022; Chen et al., 2022a; Min et al., 2022), in the adversarial case O-REPS cannot be extended to linear function



approximation without additional assumptions so algorithms are mostly based on Policy Optimization (Cai et al., 2020; Abbasi-Yadkori et al., 2019; Agarwal et al., 2020; He et al., 2022; Neu & Olkhovskaya, 2021; Luo et al., 2021; Wei et al., 2021).

**Policy Optimization in Deep RL.** Policy Optimization is among the most widely used methods in deep Reinforcement Learning (Lillicrap et al., 2015; Levine et al., 2016; Gu et al., 2017). The origins of these algorithms are Policy Gradient (Sutton & Barto, 2018), Conservative Policy Iteration (Kakade & Langford, 2002) and Natural Policy Gradient (Kakade, 2001). These have evolved into some of the state-of-the-art algorithms in RL, e.g., Trust Region Policy Optimization (TRPO) (Schulman et al., 2015), Proximal Policy Optimization (PPO) (Schulman et al., 2017) and Soft Actor-Critic (SAC) (Haarnoja et al., 2018). Recently, the connections between deep RL policy optimization algorithms and online learning regularization methods (like Follow-The-Regularized-Leader and Online-Mirror-Descent) were studied and explained (Shani et al., 2020a; Tomar et al., 2022).

**Remark A.1 (The Loop-Free Assumption)** We warn the readers that some of the works mentioned in this section (mainly in the adversarial MDP literature, e.g., Luo et al. (2021)) present a slightly different dependence in the Horizon. The reason is that they make the loop-free assumption, i.e., they assume that the state space consists of disjoint sets  $S = S_1 \cup S_2 \cup \dots \cup S_H$  such that in step  $t$  the agent can only be found in states from the set  $S_t$ . Effectively, this means that their state space is larger than ours by a factor of  $H$ . So when they present a regret bound of  $\tilde{O}(H^2 S^2 AK)$ , this implies a bound of  $\tilde{O}(H^3 S^2 AK)$  in the model presented in this paper. We emphasize that these differences are only due to different models, and not due to actual different regret bounds.

## B. Delay-Adapted Policy Optimization for (Tabular) Adversarial MDP with Known Transition

Algorithm 3 Delay-Adapted Policy Optimization with Known Transition Function (Tabular)

Input: state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , horizon  $H$ , transition function  $p$ , learning rate  $\gamma > 0$ , exploration parameter  $\epsilon > 0$ .  
 Initialization: Set  $Q_h^1(a|s) = \frac{1}{|\mathcal{A}|}$  for every  $(s; a; h) \in \mathcal{S} \times \mathcal{A} \times [H]$ .

for  $k = 1; 2; \dots; K$  do

Play episode  $k$  with policy  $\pi^k$ , observe trajectory  $(s_h^k; a_h^k)_{h=1}^H$  and compute  $Q_h^k(s)$  for every  $(s; h) \in \mathcal{S} \times [H]$ .

# Policy Evaluation

for  $j$  such that  $j + d^j = k$  do

Observe bandit feedback  $(r_h^j; a_h^j)_{h=1}^H$  and set  $B_{H+1}^j(s; a) = 0$  for every  $(s; a) \in \mathcal{S} \times \mathcal{A}$ .

for  $h = H; H-1; \dots; 1$  do

for  $(s; a) \in \mathcal{S} \times \mathcal{A}$  do

$$\text{Compute } Q_h^j(s; a) = \frac{r_h^j(s; a)}{\sum_{a'} p_h^j(s; a, a')} + \gamma \sum_{a'} p_h^j(s; a, a') Q_{h+1}^j(s; a')$$

$$\text{Compute } Q_h^j(s; a) = r_h^j(s; a) + \gamma \sum_{a'} p_h^j(s; a, a') Q_{h+1}^j(s; a')$$

$$\text{Compute } B_{h+1}^j(s; a) = Q_{h+1}^j(s; a) + \gamma \sum_{a'} p_h^j(s; a, a') B_{h+1}^j(s; a')$$

end for

end for

end for

# Policy Improvement

Define the policy  $\pi^{k+1}$  for every  $(s; a; h) \in \mathcal{S} \times \mathcal{A} \times [H]$  by:

$$\pi^{k+1}(a|s) = \frac{\exp(Q_h^k(s; a) - B_{h+1}^k(s; a))}{\sum_{a'} \exp(Q_h^k(s; a') - B_{h+1}^k(s; a'))}$$

end for

Theorem B.1. Set  $\epsilon = H^2 S A K + H^4 (K + D)^{1/2}$  and  $\gamma = 2/H$ . Running Algorithm 3 in an adversarial MDP  $M = (\mathcal{S}; \mathcal{A}; H; p; f; c^k)_{k=1}^K$  with known transition function and delayed feedback  $g_{k=1}^K$  guarantees, with probability at least  $1 - \delta$ ,

$$R_K = O\left(H^2 \frac{D}{S A K} \log \frac{K H S A}{\delta} + H^3 \frac{D}{K + D} \log \frac{K H S A}{\delta} + H^4 d_{\max} \log \frac{K H S A}{\delta}\right)$$

### B.1. The good event

Let  $\epsilon = 10 \log \frac{10 K H S A}{\delta}$ ,  $H^k$  be the history of episodes  $\{j : j + d^j < k\}$ , and define  $E_k[\cdot] = E[\cdot | H^k]$ . Define the following events:

$$E^d = \bigcap_{k=1}^K \bigcap_{j=1}^a \bigcap_{h=1}^H \left\{ \text{If } |H^k| + d^k < j + d^j \text{ then } \frac{Q_h^k(s; a) - Q_h^k(s; a)}{Q_h^k(s; a)} \leq \frac{10 H^2 d_{\max} \log \frac{10 H S A}{\delta}}{\epsilon} \right\}$$

$$E^e = \bigcap_{k=1}^K \bigcap_{h; s; a} \left\{ \frac{Q_h^k(s; a) - Q_h^k(s; a)}{Q_h^k(s; a)} \leq \frac{H^2 \log \frac{10 H S A}{\delta}}{\epsilon} \right\}$$

$$E^b = \bigcap_{k=1}^K \bigcap_{h; s; a} \left\{ \frac{Q_h^k(s; a) - Q_h^k(s; a)}{Q_h^k(s; a) + \epsilon} \leq \frac{H^2 \log \frac{10 H S A}{\delta}}{\epsilon} \right\}$$

$$E^f = \bigcap_{k=1}^K \bigcap_{h; s} \left\{ \frac{Q_h^k(s; a) - Q_h^k(s; a)}{Q_h^k(s; a)} \leq \frac{H^2 \log \frac{10 H S A}{\delta}}{\epsilon} \right\}$$

$$\frac{1}{3} \bigcap_{k=1}^K \bigcap_{h; s} \left\{ \frac{Q_h^k(s; a) - Q_h^k(s; a)}{Q_h^k(s; a)} \leq \frac{H^2 \log \frac{10 H S A}{\delta}}{\epsilon} \right\}$$

The good event is the intersection of the above events. The following lemma establishes that the good event holds with high probability.

Lemma B.2 (The Good Event) Let  $G = E^d \cap E^e \cap E^b \cap E^f$  be the good event. It holds that  $\Pr[G] \geq 1 - \epsilon$ .

Proof. We'll show that each of the events  $E^d, E^e, E^b, E^f$  holds with probability of at most  $\epsilon/4$  and so by the union bound  $\Pr[G] \geq 1 - \epsilon$ .

Event  $E^d$ : Fix  $s$  and  $h$ . For every  $(h^0, s^0, a)$  set:

$$Z_{h^0}^k(s^0, a) = \mathbb{I}\{s^0 = s; h^0 = h\} \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) r_h^k(s, a) Q_h^j(s, a)$$

$$Z_{h^0}^k(s^0, a) = \mathbb{I}\{s^0 = s; h^0 = h\} \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) r_h^k(s, a) M_h^k(s, a)$$

$$M_h^k(s, a) = \mathbb{I}\{s_h^k = s; a_h^k = a\} \sum_{h=1}^H c_h^k(s_h^k, a_h^k) + (1 - \mathbb{I}\{s_h^k = s; a_h^k = a\}) Q_h^k(s, a)$$

Note that  $\mathbb{E}_k[Z_{h^0}^k(s^0, a)] = Z_{h^0}^k(s^0, a)$  and,

$$\sum_{h^0, s^0, a} \frac{\mathbb{I}\{s_h^k = s; a_h^k = a\} \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) Q_h^j(s, a)}{c_h^k(s, a) + \sum_{h^0, s^0, a} Z_{h^0}^k(s^0, a)} = \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) Q_h^j(s, a)$$

$$\sum_{h^0, s^0, a} Z_{h^0}^k(s^0, a) \leq \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) Q_h^j(s, a);$$

where in the inequality we used the fact  $\mathbb{I}\{s_h^k = s; a_h^k = a\} \leq 1$ . Finally, we use Lemma E.5 with  $q_h^k(s, a) = c_h^k(s, a)$  and  $R = 2Hd_{\max}$  since the number of  $j$ 's such that  $j - k + d^k < j + d^j$  is at most  $2d_{\max}$ . Thus, the event holds for  $(s, h)$  with probability  $1 - \frac{\epsilon}{10HS}$ . By taking the union bound over all  $(s, h)$ ,  $E^d$  holds with probability  $1 - \frac{\epsilon}{10}$ .

Event  $E^e$  (Lemma C.2 of Luo et al. (2021)):  $E^e$  holds with probability of at least  $1 - \frac{\epsilon}{10}$  by applying Lemma E.5 with  $q_{h^0}^k(s^0, a) = c_{h^0}^k(s^0, a)$ ,  $Z_h^k(s, a) = q_h(s, a) r_h^k(s, a) Q_h^k(s, a)$  and,

$$Z_h^k(s, a) = q_h(s, a) r_h^k(s, a) \mathbb{I}\{s_h^k = s; a_h^k = a\} \sum_{h=1}^H c_{h^0}^k(s_{h^0}^k, a_{h^0}^k) + (1 - \mathbb{I}\{s_h^k = s; a_h^k = a\}) Q_h^k(s, a) :$$

Note that  $R = H, \frac{\mathbb{I}\{s_h^k = s; a_h^k = a\} \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) Q_h^j(s, a)}{c_h^k(s, a) + \sum_{h^0, s^0, a} Z_{h^0}^k(s^0, a)} = q_h(s, a) Q_h^k(s, a)$  and  $\frac{q_h^k(s, a) Z_h^k(s, a)}{c_h^k(s, a)} = q_h(s, a) Q_h^k(s, a)$ .

Event  $E^b$ : Similar to the last two events,  $E^b$  holds with probability of at least  $1 - \frac{\epsilon}{10}$  by applying Lemma E.5 with

$$q_{h^0}^k(s^0, a) = c_{h^0}^k(s^0, a) \text{ and } Z_h^k(s, a) = \frac{q_h(s) \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) r_h^k(s, a)}{c_h^k(s, a) + \sum_{h^0, s^0, a} Z_{h^0}^k(s^0, a)}.$$

Event  $E^f$ : Let  $Y_k = \sum_{h, s} \sum_{a} q_h(s) \sum_{j=1}^X \mathbb{I}\{j - k + d^k < j + d^j\} g_h^{k+d^k}(a, j, s) Q_h^k(s, a)$ . We'll use a variant of Freedman's inequality (Lemma E.3) to

bound  $\sum_{k=1}^K \mathbb{E}_k[Y_k]$ . Note that:

$$\begin{aligned} \mathbb{E}_k Y_k^2 &= \mathbb{E}_k \sum_{h,s;a} \left( q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \right)^2 \\ &= \mathbb{E}_k \sum_{h,s;a} q_h(s) \frac{r_h^k(s;a)^2}{(q_h^k(s;a))^2} \\ &= \mathbb{E}_k \sum_{h,s;a} q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \frac{r_h^k(s;a)}{q_h^k(s;a)} \\ &= \mathbb{E}_k \sum_{h,s;a} q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \frac{r_h^k(s;a)}{q_h^k(s;a)} \end{aligned}$$

where the first inequality is Cauchy-Schwartz inequality. Also  $\sum_{k=1}^K \mathbb{E}_k Y_k \leq \frac{H^2}{3} \ln \frac{10}{1-\epsilon}$ . Therefore by Lemma E.3 with probability  $1-\epsilon$ ,

$$\begin{aligned} \sum_{k=1}^K \mathbb{E}_k[Y_k] &\leq \sum_{k=1}^K \mathbb{E}_k Y_k + \frac{H^2}{3} \ln \frac{10}{1-\epsilon} \\ &= \frac{1}{3} \sum_{k=1}^K \sum_{h,s} q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} + \frac{H^2}{3} \ln \frac{10}{1-\epsilon}. \end{aligned}$$

### B.2. Proof of the main theorem

Proof of Theorem B.1 By Lemma B.2, the good event holds with probability  $1-\epsilon$ . We now analyze the regret under the assumption that the good event holds. We start with the following regret decomposition,

$$\begin{aligned} R_K &= \sum_{k=1}^K \sum_{h,s} \left( q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} - q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \right) \\ &= \underbrace{\sum_{k=1}^K \sum_{h,s} \left( q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} - q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \right)}_{\text{BIAS}_1} \\ &\quad + \underbrace{\sum_{k=1}^K \sum_{h,s} \left( q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} - q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \right)}_{\text{BIAS}_2} + \underbrace{\sum_{k=1}^K \sum_{h,s} \left( q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} - q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \right)}_{\text{BONUS}} \\ &\quad + \underbrace{\sum_{k=1}^K \sum_{h,s} \left( q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} - q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \right)}_{\text{REG}} \\ &\quad + \underbrace{\sum_{k=1}^K \sum_{h,s} \left( q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} - q_h(s) \frac{r_h^k(s;a)}{q_h^k(s;a)} \right)}_{\text{DRIFT}} \end{aligned}$$

where the first equality is by Lemma E.1 (value difference lemma),  $\mathbb{E} \sum_{k=1}^K \sum_{h=1}^H \mathbb{1}_{\text{BIAIS}_2}$  is bounded under event  $\mathbb{E}^f$  by  $O\left(\frac{H^2}{\epsilon}\right)$ . The other four terms are bounded in Lemmas B.3, B.4, B.6 and B.7. Overall,

$$\begin{aligned}
 R_K & \leq \underbrace{\frac{2}{3} \sum_{k=1}^K \sum_{h=1}^H \left( q_h(s) b_h^k(s) + O\left( H^4(K+D) + \frac{H^4 d_{\max}}{3} + \frac{H^2}{3} \right) \right)}_{\text{BIAS}_1} + O\left(\frac{H^2}{\epsilon}\right) \\
 & \quad + \underbrace{3 H^2 \text{SAK}}_{\text{BONUS}} \underbrace{\sum_{k=1}^{K-1} \sum_{h=1}^H \left( q_h(s) b_h^k(s) + O\left( H^4(K+D) + \frac{H^4 d_{\max}}{3} + \frac{H^2}{3} \right) \right)}_{\text{REG}} + \frac{H \ln A}{3} + \frac{1}{3} \sum_{k=1}^K \sum_{h=1}^H \left( q_h(s) b_h^k(s) + O\left( H^4(K+D) + \frac{H^4 d_{\max}}{3} + \frac{H^2}{3} \right) \right) \\
 & \quad + O\left( H^5(K+D) + \frac{H^5 d_{\max}}{3} \right) + O\left( \frac{H \ln A}{\epsilon} + H^2 \text{SAK} + H^5(K+D) + \frac{H^3}{2} + \frac{H^2}{\epsilon} + \frac{H^5 d_{\max}}{\epsilon} \right) :
 \end{aligned}$$

For  $\epsilon = \frac{1}{H^2 \text{SAK} + H^4(K+D)}$  and  $\epsilon = 2H$ , we get:  $R_K \leq O\left( H^2 \overline{\text{SAK}} + H^3 \overline{D} + H^3 \overline{K} + H^4 d_{\max} \right) : \quad \square$

### B.3. Bound on $\text{BIAS}_1$

Lemma B.3. Under the good event  $\mathbb{E}^f$ ,  $\mathbb{E} \sum_{k=1}^K \sum_{h=1}^H \left( q_h(s) b_h^k(s) + O\left( H^4(K+D) + \frac{H^4 d_{\max}}{3} + \frac{H^2}{3} \right) \right) : \quad \square$

Proof. Let  $Y_k = \sum_{h=1}^H q_h(s) \sum_{j=1}^{k+d^k} \left( \mathbb{1}_{\text{BIAIS}_1} \right) \mathbb{1}_{\text{BIAIS}_1} \left( j, s; Q_h^k(s) \right)$ . It holds that

$$\text{BIAS}_1 = \sum_{k=1}^K \sum_{h=1}^H \left( q_h(s) \sum_{j=1}^{k+d^k} \left( \mathbb{1}_{\text{BIAIS}_1} \right) \mathbb{1}_{\text{BIAIS}_1} \left( j, s; Q_h^k(s) \right) \right) - \sum_{k=1}^K \mathbb{E}_k[Y_k] + \sum_{k=1}^K \mathbb{E}_k[Y_k] - \sum_{k=1}^K Y_k :$$

Under event  $\mathbb{E}^f$  it holds that

$$\sum_{k=1}^K \mathbb{E}_k[Y_k] - \sum_{k=1}^K Y_k \leq \frac{1}{3} \sum_{k=1}^K \sum_{h=1}^H \left( q_h(s) b_h^k(s) + \frac{H^2}{3} \ln \frac{10}{\epsilon} \right) :$$

In addition,

$$\begin{aligned}
 & \sum_{k=1}^K \sum_{h;s} \mathbb{E} \left[ \sum_{k=1}^D \mathbf{q}_h^k(s) \left( \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) - 1 \right) \right] \mathbb{E}_k[Y_k] = \\
 &= \sum_{k=1}^K \sum_{h;s;a} \mathbf{q}_h^k(s) \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) \left( 1 - \frac{\mathbf{q}_h^k(s; a) r_h^k(s; a)}{\mathbf{q}_h^k(s; a) + \dots} \right) \\
 &= \sum_{k=1}^K \sum_{h;s;a} \mathbf{q}_h^k(s) \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) \left( 1 - \frac{(\mathbf{q}_h^k(s; a) + \dots) r_h^k(s; a)}{\mathbf{q}_h^k(s; a) + \dots} \right) \\
 &\quad + \sum_{k=1}^K \sum_{h;s;a} \mathbf{q}_h^k(s) \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) \frac{r_h^k(s; a)}{\mathbf{q}_h^k(s; a) + \dots} \\
 &= \sum_{k=1}^K \sum_{h;s;a} \mathbf{q}_h^k(s) \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) (1 - r_h^k(s; a)) + \sum_{k=1}^K \sum_{h;s;a} \mathbf{q}_h^k(s) \frac{H \sum_{j \in \mathcal{A}} r_h^k(s; a)}{\mathbf{q}_h^k(s; a) + \dots} \\
 &= H \sum_{k=1}^K \sum_{h;s;a} \mathbf{q}_h^k(s) \max_{j \in \mathcal{A}} \left( \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) - \sum_{j \in \mathcal{A}} r_h^k(s; a) \right) + \sum_{k=1}^K \sum_{h;s;a} \mathbf{q}_h^k(s) \frac{H \sum_{j \in \mathcal{A}} r_h^k(s; a)}{\mathbf{q}_h^k(s; a) + \dots} \\
 &= H \sum_{k=1}^K \sum_{h;s} \mathbf{q}_h^k(s) \left( \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) - \sum_{j \in \mathcal{A}} r_h^k(s; a) \right) + \frac{1}{3} \sum_{k=1}^K \sum_{h;s} \mathbf{q}_h^k(s) b_h^k(s).
 \end{aligned}$$

Finally, using Lemma B.8,

$$\begin{aligned}
 \text{BIAS}_1 &= \frac{2}{3} \sum_{k=1}^K \sum_{h;s} \mathbf{q}_h^k(s) b_h^k(s) + H \sum_{h;s} \mathbf{q}_h^k(s) \sum_{k=1}^K \left( \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) - \sum_{j \in \mathcal{A}} r_h^k(s; a) \right) + \frac{H^2}{\dots} \\
 &= \frac{2}{3} \sum_{k=1}^K \sum_{h;s} \mathbf{q}_h^k(s) b_h^k(s) + H \sum_{h;s} \mathbf{q}_h^k(s) \mathcal{O} \left( H^2(K+D) + \frac{H^2 d_{\max}}{\dots} \right) + \frac{H^2}{\dots} \\
 &= \frac{2}{3} \sum_{k=1}^K \sum_{h;s} \mathbf{q}_h^k(s) b_h^k(s) + \mathcal{O} \left( H^4(K+D) + \frac{H^4 d_{\max}}{\dots} \right) + \frac{H^2}{\dots}. \quad \square
 \end{aligned}$$

#### B.4. Bound on $\text{BONUS}$

Lemma B.4. It holds that  $\text{BONUS} \leq 3H^2 \text{SAK} \sum_{k,h;s} \mathbf{q}_h^k(s) b_h^k(s)$ :

Proof. Note that  $\mathbf{B}_h^k$  is the Q-function of policy  $\pi^k$  with respect to the cost function  $c^k$ . Hence, by the value difference lemma (Lemma E.1),

$$\sum_{h;s} \mathbf{q}_h^k(s) b_h^k(s) = \sum_{h;s} \mathbf{q}_h^k(s) \left( V_1^k(s_{\text{init}}; b^k) - V(s_{\text{init}}; b^k) \right) = \sum_{h;s} \mathbf{q}_h^k(s) \left( \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) - \sum_{j \in \mathcal{A}} r_h^k(s; a) \right) :$$

Summing over  $k$  we get:  $\text{BONUS} = \sum_{k,h;s} \mathbf{q}_h^k(s) b_h^k(s) = \sum_{k,h;s} \mathbf{q}_h^k(s) \left( \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) - \sum_{j \in \mathcal{A}} r_h^k(s; a) \right)$ : For last,

$$\sum_{k=1}^K \sum_{h;s} \mathbf{q}_h^k(s) b_h^k(s) = 3H \sum_{k=1}^K \sum_{h;s;a} \frac{\mathbf{q}_h^k(s) \sum_{j \in \mathcal{A}} r_h^k(s; a)}{\mathbf{q}_h^k(s) \sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) + \dots} \leq 3H^2 \text{SAK} :$$

where the last uses the fact that  $\frac{\sum_{j \in \mathcal{A}} r_h^k(s; a)}{\sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) + \dots} \leq \sum_{j \in \mathcal{A}} r_h^k(s; a)$ .  $\square$

Remark B.5. The adaptation to delay via the ratio  $\frac{\sum_{j \in \mathcal{A}} r_h^k(s; a)}{\sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) + \dots}$  is simple, yet crucial. The main reason is the following. While in MAB the ratio  $\frac{\sum_{j \in \mathcal{A}} r_h^k(s; a)}{\sum_{j \in \mathcal{A}} \mathbf{Q}_h^k(s; a) + \dots} = \sum_{j \in \mathcal{A}} r_h^k(s; a)$  is always bounded by a constant (Thune et al., 2019, Lemma 11), in MDPs it can be

as large as  $e^{d_{\max}}$ . In fact, even the ratio  $\frac{Q_h^{k+1}(a_j|s)}{Q_h^k(a_j|s)} = \frac{Q_h^k(a_j|s)}{Q_h^k(a_j|s)}$  can be of order  $e^{d_{\max}}$ , because even if an action is chosen with probability close to 1, the estimator  $Q_h^k(s; a)$  can be as large as  $e^{d_{\max}}$ , as long as the visitation probability to states is smaller than  $\frac{1}{e^{d_{\max}}}$ . This can cause radical changes in the probability to take an action and as a consequence in the probability of the rest of the actions in that state). For example, assume we have two actions with  $Q_h^k(a_1|s) = 1 = (e^{d_{\max}} + 1)^{-1}$ ,  $Q_h^k(a_2|s) = e^{-d_{\max}} = (e^{d_{\max}} + 1)^{-1}$  and  $Q_h^k(s)$ . Now, assume that the feedback from  $d_{\max}$  episodes arrive at the end of episode  $k$  for which the agent visited in the state-action pair  $(s, a_1)$ . Further assume the cost-to-go from  $(s, a_1)$  was of order  $H$ . This would imply that  $Q_h^{k+1}(a_1|s) = Q_h^k(a_1|s) + H = (e^{d_{\max}} + 1)^{-1} + H = (e^{d_{\max}} + 1)^{-1} + H$ . In particular,  $\frac{Q_h^{k+1}(a_2|s)}{Q_h^k(a_2|s)} = \frac{(e^{d_{\max}} + 1)^{-1} + H}{(e^{d_{\max}} + 1)^{-1}} = (e^{d_{\max}} + 1) + H$ .

B.5. Bound on REG

Lemma B.6. For  $\frac{H \ln A}{2H}$  it holds that  $\text{REG} \leq \frac{H \ln A}{3} + \frac{1}{3} \sum_{k,h;s} q_h(s) b_h^k(s) + O(H^5 K) + \frac{H^3}{2}$  :

Proof. By Corollary E.7, since  $\max_{k,h;s,a} B_h^k(s; a) \leq 3H^2$ ,

$$\begin{aligned} \text{REG} &\leq \frac{H \ln A}{2} + 2 \sum_{k,h;s,a} q_h(s) \frac{1}{h^{k+d^k}} (a_j|s) Q_h^k(s; a) B_h^k(s; a)^2 + O(H^5 K) \\ &\leq \frac{H \ln A}{2} + 2 \sum_{k,h;s,a} q_h(s) \frac{1}{h^{k+d^k}} (a_j|s) Q_h^k(s; a)^2 + 2 \sum_{k,h;s,a} q_h(s) \frac{1}{h^{k+d^k}} (a_j|s) B_h^k(s; a)^2 + O(H^5 K): \quad (17) \end{aligned}$$

For the middle term

$$\begin{aligned} 2 \sum_{k,h;s,a} q_h(s) \frac{1}{h^{k+d^k}} (a_j|s) Q_h^k(s; a)^2 &\leq 2 \sum_{k,h;s,a} q_h(s) \frac{1}{h^{k+d^k}} (a_j|s) \frac{H^2 r_h^k(s; a)^2 \mathbb{1}\{s_h^k = s; a_h^k = a\}}{(Q_h^k(s; a) + \frac{1}{h})^2} \\ &\leq 2 H^2 \sum_{k,h;s,a} \frac{q_h(s) \frac{1}{h^{k+d^k}} (a_j|s) r_h^k(s; a)}{Q_h^k(s; a) + \frac{1}{h}} + O\left(\frac{H^3}{2}\right) \\ &\leq \frac{2}{3} H \sum_{k,h;s} q_h(s) b_h^k(s) + O\left(\frac{H^3}{2}\right) \\ &\leq \frac{1}{3} \sum_{k,h;s} q_h(s) b_h^k(s) + O\left(\frac{H^3}{2}\right); \end{aligned}$$

where the second inequality is by  $\frac{1}{h}$  and the last is since  $\frac{H \ln A}{2H}$ . For the last term in Eq. (17) we use  $\max_{k,h;s} B_h^k(s) \leq 3H$  and therefore  $\sum_{k,h;s,a} q_h(s) \frac{1}{h^{k+d^k}} (a_j|s) B_h^k(s; a)^2 \leq 9 H^5 K$ . Overall,

$$\text{REG} \leq \frac{H \ln A}{3} + \frac{1}{3} \sum_{k,h;s} q_h(s) b_h^k(s) + O(H^5 K) + \frac{H^3}{2} : \quad \square$$

B.6. Bound on DRIFT

Lemma B.7. If event  $E^d$  holds then  $\text{DRIFT} \leq O(H^5(K + D)) + \frac{H^5 d_{\max}}{2}$  :

Proof. We use Lemma B.8 and the fact that  $\|Q_h^k(s; a) - B_h^k(s; a)\| \leq 3H^2$  to obtain

$$\begin{aligned} \text{DRIFT} &= \sum_{k=1}^K \sum_{h;s;a} Q_h(s) \left( \sum_{j \in \mathcal{A}} k_h^k(a; j; s) - \sum_{j \in \mathcal{A}} k_h^{k+d^k}(a; j; s) \right) (Q_h^k(s; a) - B_h^k(s; a)) \\ &\quad - \sum_{k=1}^K \sum_{h;s;a} Q_h(s) \sum_{j \in \mathcal{A}} \left( k_h^k(a; j; s) - \sum_{j \in \mathcal{A}} k_h^{k+d^k}(a; j; s) \right) (Q_h^k(s; a) - B_h^k(s; a)) \\ &\leq 3H^2 \sum_{k=1}^K \sum_{h;s} Q_h(s) \sum_{j \in \mathcal{A}} \left( k_h^k(j; s) - \sum_{j \in \mathcal{A}} k_h^{k+d^k}(j; s) \right) k_1 \\ &\leq H^5(K + D) + \frac{H^5 d_{\max}}{2} \quad \square \end{aligned}$$

Lemma B.8. If event  $E^d$  holds then,  $\sum_{k=1}^K \sum_{h;s} \left( k_h^{k+d^k}(j; s) - \sum_{j \in \mathcal{A}} k_h^k(j; s) \right) k_1 \leq H^2(K + D) + \frac{H^2 d_{\max}}{2}$  :

Proof. We first bound,

$$\begin{aligned} \sum_{k=1}^K \sum_{h;s} \left( k_h^{k+d^k}(j; s) - \sum_{j \in \mathcal{A}} k_h^k(j; s) \right) k_1 &= \sum_{k=1}^K \sum_{j=k}^{k+d^k-1} \sum_{h;s} \left( k_h^{j+1}(j; s) - \sum_{j \in \mathcal{A}} k_h^j(j; s) \right) k_1 \\ &= \sum_{k=1}^K \sum_{j=k}^{k+d^k-1} \sum_{h;s} \sum_{a \in \mathcal{A}} \left( j_h^{j+1}(a; j; s) - \sum_{j \in \mathcal{A}} j_h^j(a; j; s) \right) \end{aligned}$$

Now, we apply Lemma E.9 for each  $j$  in the summation above with  $\tilde{Q}_h^j(\cdot) = j_h^{j+1}(\cdot; j; s)$ ,  $\tilde{B}_h^j(\cdot) = \sum_{j \in \mathcal{A}} j_h^j(\cdot; j; s)$  and  $\tilde{C}_h^j(\cdot) = \sum_{i:i+d^i=j} (Q_h^i(s; \cdot) - B_h^i(s; \cdot))$ . We observe that,

$$\begin{aligned} \sum_{k=1}^K \sum_{h;s} \left( k_h^{k+d^k}(j; s) - \sum_{j \in \mathcal{A}} k_h^k(j; s) \right) k_1 &= \sum_{k=1}^K \sum_{j=k}^{k+d^k-1} \sum_{h;s} \sum_{a \in \mathcal{A}} \left( j_h^{j+1}(a; j; s) - \sum_{j \in \mathcal{A}} j_h^j(a; j; s) \right) \sum_{i:i+d^i=j} (Q_h^i(s; a) + 6H^2) \\ &\quad + \sum_{k=1}^K \sum_{j=k}^{k+d^k-1} \sum_{h;s} \sum_{a \in \mathcal{A}} \left( j_h^{j+1}(a; j; s) - \sum_{j \in \mathcal{A}} j_h^j(a; j; s) \right) \sum_{i:i+d^i=j} (Q_h^i(s; a^0) + 6H^2) \\ &= 2 \sum_{k=1}^K \sum_{j=k}^{k+d^k-1} \sum_{h;s} \sum_{a \in \mathcal{A}} \left( j_h^{j+1}(a; j; s) - \sum_{j \in \mathcal{A}} j_h^j(a; j; s) \right) \sum_{i:i+d^i=j} (Q_h^i(a; j; s) + 6H^2) \\ &= 2 \sum_{k=1}^K \sum_{j=k}^{k+d^k-1} \sum_{i:i+d^i=j} \sum_{a \in \mathcal{A}} \left( j_h^{j+1}(a; j; s) - \sum_{j \in \mathcal{A}} j_h^j(a; j; s) \right) \sum_{i:i+d^i=j} (Q_h^i(a; j; s) + 12H^2) \quad \text{If } k + i + d^i < k + d^k \\ &\leq 2 \sum_{k=1}^K \sum_{j=k}^{k+d^k-1} \sum_{i:i+d^i=j} \sum_{a \in \mathcal{A}} \left( j_h^{j+1}(a; j; s) - \sum_{j \in \mathcal{A}} j_h^j(a; j; s) \right) \sum_{i:i+d^i=j} (Q_h^i(a; j; s) + 12H^2(D + K)); \end{aligned}$$



where the last inequality is by Lemma E.10. For the first term we use  $\text{Event}$

$$\begin{aligned}
 \sum_{k=1}^K \sum_{j=k}^K \sum_{i:i+d^i=j}^K \sum_a \mathbb{E} \left[ \sum_{h=1}^H (a_j^i(s) Q_h^i(a_j^i(s)) - \sum_{k=1}^K \sum_{j=k}^K \sum_{i:i+d^i=j}^K \sum_a \mathbb{E} \left[ \sum_{h=1}^H (a_j^i(s) Q_h^i(a_j^i(s)) \right] \right) \right] \\
 = \sum_{k=1}^K \sum_{i=1}^K \sum_a \mathbb{E} \left[ \sum_{h=1}^H (a_j^i(s) Q_h^i(a_j^i(s)) - \sum_{k=1}^K \sum_{j=k}^K \sum_{i:i+d^i=j}^K \sum_a \mathbb{E} \left[ \sum_{h=1}^H (a_j^i(s) Q_h^i(a_j^i(s)) \right] \right) \right] + O \left( \frac{H^2 d_{\max}}{K} \right) \\
 \leq H \sum_{k=1}^K \sum_{i=1}^K \mathbb{E} \left[ \sum_{h=1}^H (a_j^i(s) Q_h^i(a_j^i(s)) - \sum_{k=1}^K \sum_{j=k}^K \sum_{i:i+d^i=j}^K \sum_a \mathbb{E} \left[ \sum_{h=1}^H (a_j^i(s) Q_h^i(a_j^i(s)) \right] \right) \right] + O \left( \frac{H^2 d_{\max}}{K} \right) \\
 \leq H(D + K) + O \left( \frac{H^2 d_{\max}}{K} \right) : \quad \square
 \end{aligned}$$

## C. Delay-Adapted Policy Optimization for (Tabular) Adversarial MDP with Unknown Transition

## Algorithm 4 Delay-Adapted Policy Optimization with Unknown Transition Function (Tabular)

Input: state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , horizon  $H$ , learning rate  $\gamma > 0$ , exploration parameter  $\epsilon > 0$ , confidence parameter  $\delta > 0$ .

Initialization: Set  $q_h^1(a|s) = \frac{1}{|\mathcal{A}|}$  for every  $(s; a; h) \in \mathcal{S} \times \mathcal{A} \times [H]$ .

for  $k = 1; 2; \dots; K$  do

Play episode  $k$  with policy  $q_h^k$  and observe trajectory  $(s_h^k; a_h^k)_{h=1}^H$ .

Compute visit counters for every  $(h; s; a; s^0) \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ :

$$n_h^k(s; a; s^0) = \sum_{j: j+d^j < k} \mathbb{1}\{s_h^j = s; a_h^j = a; s_{h+1}^j = s^0\} g_{h=1}^H; \quad n_h^k(s; a) = \sum_{j: j+d^j < k} \mathbb{1}\{s_h^j = s; a_h^j = a\} g_{h=1}^H$$

Compute empirical transition function  $p_h^k(s^0|s; a) = \frac{n_h^k(s; a; s^0)}{\max\{n_h^k(s; a); 1\}g_{h=1}^H}$  and confidence set  $\mathcal{P}_h^k = \{p_h^k(s; a)g_{s; a; h} \text{ such that } p_h^0(s^0|s; a) \in \mathcal{P}_h^k(s; a) \text{ if and only if } \sum_{s^0} p_h^0(s^0|s; a) = 1 \text{ and for every } s^0 \in \mathcal{S}\}$ :

$$\mathcal{P}_h^k(s^0|s; a) = \mathcal{P}_h^k(s^0|s; a) \pm 4 \frac{s}{n_h^k(s; a) - 1} \log \frac{10HSAK}{n_h^k(s; a) - 1} + \frac{10 \log \frac{10HSAK}{n_h^k(s; a) - 1}}{n_h^k(s; a) - 1}$$

Compute occupancy measures  $\mu_h^k(s) = \max_{p \in \mathcal{P}_h^k} q_h^k \cdot p^0(s)$  and  $\underline{\mu}_h^k(s) = \min_{p \in \mathcal{P}_h^k} q_h^k \cdot p^0(s)$ .

# Policy Evaluation

for  $j$  such that  $j + d^j = k$  do

Observe bandit feedback  $(s_h^j; a_h^j)_{h=1}^H$  and set  $B_{H+1}^j(s; a) = 0$  for every  $(s; a) \in \mathcal{S} \times \mathcal{A}$ .

for  $h = H; H-1; \dots; 1$  do

for  $(s; a) \in \mathcal{S} \times \mathcal{A}$  do

Compute  $q_h^j(s; a) = \frac{q_h^j(a|s)}{\max\{q_h^j(a|s); \frac{\delta}{k(a|s)g_{h=1}^H}\}}$  and  $L_h^j = \prod_{h^0=h}^H q_{h^0}^j(s_{h^0}^j; a_{h^0}^j)$ .

Compute  $\bar{q}_h^j(s) = \prod_{a \in \mathcal{A}} \frac{\sum_{h^0=h}^H q_{h^0}^j(a|s) r_{h^0}^j(s; a) (\bar{q}_{h^0}^j(s; a) - q_{h^0}^j(s; a))}{\sum_{h^0=h}^H q_{h^0}^j(a|s) + \frac{\delta}{k(a|s)g_{h^0=1}^H}}$  and  $\underline{q}_h^j(s) = \prod_{a \in \mathcal{A}} \frac{\sum_{h^0=h}^H q_{h^0}^j(a|s) r_{h^0}^j(s; a) (\underline{q}_{h^0}^j(s; a) - q_{h^0}^j(s; a))}{\sum_{h^0=h}^H q_{h^0}^j(a|s) + \frac{\delta}{k(a|s)g_{h^0=1}^H}}$ .

Compute  $\bar{b}_h^j(s) = \bar{q}_h^j(s) + \underline{q}_h^j(s)$  and  $\bar{Q}_h^j(s; a) = r_h^j(s; a) \frac{\bar{q}_h^j(s)}{\sum_{h^0=h}^H q_{h^0}^j(a|s) + \frac{\delta}{k(a|s)g_{h^0=1}^H}}$ .

Compute  $B_h^j(s; a) = \bar{b}_h^j(s) + \max_{p \in \mathcal{P}_h^k} \sum_{s^0 \in \mathcal{S}} p_h^0(s^0|s; a) \sum_{a^0 \in \mathcal{A}} B_{h+1}^j(s^0, a^0)$ .

end for

end for

end for

# Policy Improvement

Define the policy  $q_h^{k+1}$  for every  $(s; a; h) \in \mathcal{S} \times \mathcal{A} \times [H]$  by:

$$q_h^{k+1}(a|s) = \frac{\prod_{j: j+d^j=k} \exp(\bar{Q}_h^j(s; a) - B_h^j(s; a))}{\prod_{j: j+d^j=k} \left( \sum_{a^0 \in \mathcal{A}} \exp(\bar{Q}_h^j(s; a^0) - B_h^j(s; a^0)) \right)}$$

end for

Theorem C.1. Set  $\epsilon = H^{-2} \sqrt{SAK} + H^4(K+D)^{1/2}$  and  $\delta = 2H$ . Running Algorithm 4 in an adversarial MDP  $M = (\mathcal{S}; \mathcal{A}; H; p; f; c; g_{k=1}^K)$  with unknown transition function and delayed feedback  $g_{k=1}^K$  guarantees, with probability at least  $1 - \epsilon$ ,

$$R_K \leq O \left( H^3 S^D \frac{1}{AK} \log \frac{SAK}{\epsilon} + H^3 D \frac{1}{K+D} \log \frac{SAK}{\epsilon} + H^4 S^2 A d_{\max} + H^4 S^3 A \log^2 \frac{SAK}{\epsilon} \right)$$

Remark C.2 (Dependence on  $d_{\max}$ ). All our regret bounds contain additive terms that scale linearly with  $d_{\max}$ . While these are low-order terms when  $d_{\max}$  is smaller than  $D$ , they may become dominant for large maximal delay. The dependence on  $d_{\max}$  can be removed altogether using the clipping technique (Thune et al., 2019; Bistriz et al., 2019; Lancewicki

et al., 2022b), i.e., ignoring episodes that their delay is larger than some threshold. In the case of known transitions (Theorem B.1 in Appendix B), we can set  $\bar{D} = H$  and remove the dependence on  $d_{\max}$  without hurting our original regret bound, i.e., we get the bound  $R_K = \mathcal{O}(H^2 S A K + H^3 S^2 (K + D))$ . However, in the case of unknown transitions (Theorem C.1 in Appendix C), we get a slightly worse regret bound. Specifically, we can set  $\frac{D}{H^2 S^2 A}$  and obtain the regret  $R_K = \mathcal{O}(H^3 S^2 A (K + D))$ . Jin et al. (2022) encounter the same issue in their regret bounds, so it remains an open problem whether the dependence on  $d_{\max}$  can be removed without hurting the original regret bound in the unknown transitions case.

C.1. The good event

Let  $\beta = 10 \log \frac{10KHSA}{\epsilon}$ ,  $H^k$  be the history of episodes  $\{j : j + d^j < k\}$ , and  $H^k$  be the history of episodes  $\{j : j < k\}$ . Define the following events:

$$\begin{aligned}
 E^p &= \bigcap_{k=1}^K \bigcap_{h,s,a} \left\{ p_h(s^0; j; s; a) - p_h^k(s^0; j; s; a) \leq \frac{S}{4} \frac{p_h^k(s^0; j; s; a) \log \frac{10HSAK}{\epsilon}}{\max_h n_h^k(s; a); 1g} + 10 \frac{\log \frac{10HSAK}{\epsilon}}{\max_h n_h^k(s; a); 1g} \right\} \\
 E^{est} &= \bigcap_{k=1}^K \bigcap_{h,s,a} \left\{ |q_h^k(s; a) - \bar{q}_h^k(s; a)| \leq \frac{r}{H^4 S^2 A K \log \frac{10KHSA}{\epsilon}} + H^3 S^3 A \log^2 \frac{10KHSA}{\epsilon} + H^3 S^2 A d_{\max} \right\} \\
 E^d &= \bigcap_{k=1}^K \bigcap_{h,s} \left\{ \prod_{j=1}^k \prod_{a} \mathbb{1}_{\{j + d^j < k + d^k\}} \prod_{h} g_h^{k+d^k}(a; j; s) \left( \bar{Q}_h^k(s; a) - Q_h^k(s; a) \right) \leq \frac{10H^2 d_{\max} \log \frac{10HS}{\epsilon}}{\epsilon} \right\} \\
 E &= \bigcap_{k=1}^K \bigcap_{h,s,a} \left\{ q_h(s; a) - \bar{Q}_h^k(s; a) - Q_h^k(s; a) \leq \frac{H^2 \log \frac{10HSA}{\epsilon}}{\epsilon} \right\} \\
 E^b &= \bigcap_{k=1}^K \bigcap_{h,s,a} \left\{ \frac{q_h(s) \prod_{h}^{k+d^k}(a; j; s) r_h^k(s; a) \mathbb{1}_{\{s_h^k = s; a_h^k = a\}}}{(q_h^k(s; a) + \epsilon)^2} \leq \prod_{k=1}^k \bigcap_{h,s,a} \frac{q_h(s) \prod_{h}^{k+d^k}(a; j; s) r_h^k(s; a)}{q_h^k(s; a) + \epsilon} \leq \frac{H}{2} \ln \frac{10H}{\epsilon} \right\} \\
 E^f &= \bigcap_{k=1}^K \bigcap_{h,s} \left\{ E_k \prod_{h,s} \frac{D}{4} q_h(s) \prod_{h}^{k+d^k}(j; s); \bar{Q}_h^k(s; a) \leq \prod_{k=1}^k \bigcap_{h,s} \frac{D}{5} q_h(s) \prod_{h}^{k+d^k}(j; s); \bar{Q}_h^k(s; a) \leq \prod_{k=1}^k \bigcap_{h,s} \frac{D}{3} q_h(s) \prod_{h}^{k+d^k}(j; s) + \frac{H^2}{\epsilon} \ln \frac{10}{\epsilon} \right\}
 \end{aligned}$$

The good event is the intersection of the above events. The following lemma establishes that the good event holds with high probability.

**Lemma C.3 (The Good Event)** Let  $G = E^p \setminus E^{est} \setminus E^d \setminus E \setminus E^b \setminus E^f$  be the good event. It holds that  $\Pr[G] \geq 1 - \epsilon$ . Moreover, under the good event, it holds that  $p^k$  and  $q_h^k(s; a) - \bar{q}_h^k(s; a) - \bar{q}_h^k(s; a)$  for every  $(k; h; s; a) \in [K] \times [H] \times S \times A$ .

**Proof.** We'll show that each of the events  $E^p; E^d; E; E^b; E^f$  holds with probability of at most  $\epsilon/5$  and so by the union bound  $\Pr[G] \geq 1 - \epsilon$ .

**Event  $E^p$ :** Holds with probability  $1 - \epsilon/5$  by standard Bernstein inequality (see, e.g., Lemma 2 in Jin et al. (2020a)). As a consequence of event  $E^p$ ,  $p^k \geq 2P^k$  for all  $k$ . In particular,  $q_h^k(s; a) - \bar{q}_h^k(s; a) - \bar{q}_h^k(s; a)$  for all  $k; h; s$  and  $a$ .

**Event  $E^{est}$ :** Holds with probability  $1 - \epsilon/5$  by Jin et al. (2022, Lemma D.12) (see also (Jin et al., 2020a, Lemma 4)) which is a standard techniques (adapted to delays) of summing the the confidence radius on the trajectory.

**Event  $E^d$ :** We show the proof under event  $E^p$  which occurs with probability  $1 - \epsilon/5$ . Fix  $s$  and  $h$ . Similar to the proof of

Lemma B.2, we apply Lemma E.5. For every  $(h^0, s^0, a)$ , set  $q_{h^0}^k(s^0; a) = q_{h^0}^k(s^0; a)$  and

$$\begin{aligned} z_{h^0}^k(s^0; a) &= \text{If } s^0 = s; h^0 = hg \sum_{j=1}^X \text{If } j \quad k + d^k < j + d^j g \quad \frac{r_h^k(s; a) Q_h^j(s; a)}{q_h^k(s; a) +} \\ Z_{h^0}^k(s^0; a) &= \text{If } s^0 = s; h^0 = hg \sum_{j=1}^X \text{If } j \quad k + d^k < j + d^j g \quad \frac{r_h^k(s; a) M_h^k(s; a)}{q_h^k(s; a) +} \\ M_h^k(s; a) &= \text{If } s_h^k = s; a_h^k = ag \sum_{h=1}^X c_h^k(s_h^k; a_h^k) + (1 \quad \text{If } s_h^k = s; a_h^k = ag) Q_h^k(s; a): \end{aligned}$$

Note that  $\mathbb{E}_k[Z_{h^0}^k(s^0; a)] = z_{h^0}^k(s^0; a)$  and

$$\begin{aligned} \sum_{h^0; s^0; a} \frac{\text{If } s_h^k = s; a_h^k = sg Z_{h^0}^k(s; a)}{q_h^k(s; a) +} &= \sum_{j=1}^X \sum_a \text{If } j \quad k + d^k < j + d^j g \quad \frac{r_h^k(s; a) Q_h^j(s; a)}{q_h^k(s; a) +} \\ \sum_{h^0; s^0; a} \frac{c_h^k(s; a) z_{h^0}^k(s; a)}{q_h^k(s; a)} &= \sum_{j=1}^X \sum_a \text{If } j \quad k + d^k < j + d^j g \quad \frac{r_h^k(s; a) Q_h^k(s; a)}{q_h^k(s; a)}; \end{aligned}$$

where in the inequality we used the fact that  $q_h^k(s; a) = q_h^k(s; a)$  under the event  $\mathbb{E}^P$ . Finally, we use Lemma E.5 with  $R = 2Hd_{\max}$  as in the proof of Lemma B.2. Thus, the event holds with probability  $1 - \frac{1}{10HS}$ . By taking the union bound over all  $h$  and  $s$ ,  $\mathbb{E}^d$  holds with probability  $1 - \frac{1}{10}$ .

Event  $\mathbb{E}^a$  (Lemma C.2 of Luo et al. (2021)): We show the proof under event  $\mathbb{E}^P$  which occurs with probability  $1 - 10$ . Again,  $\mathbb{E}^a$  holds with probability of at least  $1 - 10$  by applying Lemma E.5 with  $q_{h^0}^k(s^0; a) = q_{h^0}^k(s^0; a)$ ,  $z_{h^0}^k(s; a) = q_h(s; a) r_h^k(s; a) Q_h^k(s; a)$  and,

$$Z_{h^0}^k(s; a) = q_h(s; a) r_h^k(s; a) \text{If } s_h^k = s; a_h^k = ag \sum_{h^0=1}^X c_{h^0}^k(s_{h^0}^k; a_{h^0}^k) + (1 \quad \text{If } s_h^k = s; a_h^k = ag) Q_h^k(s; a) :$$

Note that  $R = H$ ,  $\frac{\text{If } s_h^k = s; a_h^k = sg Z_{h^0}^k(s; a)}{q_h^k(s; a) +} = q_h(s; a) Q_h^k(s; a)$  and  $\frac{c_h^k(s; a) z_{h^0}^k(s; a)}{q_h^k(s; a)} = q_h(s; a) Q_h^k(s; a)$  where similar to before, the inequality holds under event  $\mathbb{E}^P$ .

Event  $\mathbb{E}^b$ : We show the proof under event  $\mathbb{E}^P$  which occurs with probability  $1 - 10$ . Similar to before,  $\mathbb{E}^b$  holds with probability of at least  $1 - 10$  by applying Lemma E.5 with  $q_{h^0}^k(s^0; a) = q_{h^0}^k(s^0; a)$  and  $z_{h^0}^k(s; a) = \frac{q_h(s) r_h^{k+d^k}(a; j; s) r_h^k(s; a)}{q_h^k(s; a) +}$ .

Event  $\mathbb{E}^f$ : We show the proof under event  $\mathbb{E}^P$  which occurs with probability  $1 - 10$ . Let  $Y_k = \sum_{h; s} q_h(s) r_h^{k+d^k}(a; j; s) Q_h^k(s; a)$ . Similar to the proof of Lemma B.2, we use a variant of Freedman's inequality (Lemma E.3) to bound  $\sum_{k=1}^K \mathbb{E}_k[Y_k]$ .

$$\begin{aligned} \mathbb{E}_k Y_k^2 &= \mathbb{E}_k \sum_{h; s; a} q_h(s) r_h^{k+d^k}(a; j; s) Q_h^k(s; a) A \sum_{h; s; a} q_h(s) r_h^{k+d^k}(a; j; s) Q_h^k(s; a) A \\ &= \mathbb{E}_k \sum_{h; s; a} q_h(s) r_h^{k+d^k}(a; j; s) A \sum_{h; s; a} q_h(s) r_h^{k+d^k}(a; j; s) (Q_h^k(s; a))^2 A \\ &= H \mathbb{E}_k \sum_{h; s; a} q_h(s) r_h^{k+d^k}(a; j; s) \frac{r_h^k(s; a)^2 (L_h^k)^2 \text{If } s_h^k = s; a_h^k = ag}{(q_h^k(s; a) +)^2} \\ &= H^3 \sum_{h; s; a} q_h(s) r_h^{k+d^k}(a; j; s) \frac{r_h^k(s; a) q_h^k(s; a)}{(q_h^k(s; a) +)^2} = H^3 \sum_{h; s; a} \frac{q_h(s) r_h^{k+d^k}(a; j; s) r_h^k(s; a)}{q_h^k(s; a) +}; \end{aligned}$$

where the first inequality is Cauchy-Schwartz inequality and the last inequality holds under  $\mathcal{E}$ . Also,  $\sum_{k=1}^K \sum_{h;s} |Y_k| \leq \frac{H^2}{3}$ . Therefore by Lemma E.3 with probability  $1 - \frac{1}{10}$ ,

$$\sum_{k=1}^K \mathbb{E}_k[Y_k] \leq \sum_{k=1}^K \sum_{h;s} \frac{H \sum_{j=1}^D q_h(s) \frac{h^{k+d^k} (a_j s) r_h^k(s; a)}{q_h^k(s; a)} A + \frac{H^2 \ln 10}{3}}{3} = \frac{1}{3} \sum_{k=1}^K \sum_{h;s} q_h(s) b_h^k(s) + \frac{H^2 \ln 10}{3}$$

## C.2. Proof of the main theorem

Proof of Theorem C.1 By Lemma C.3, the good event holds with probability of at least  $1 - \frac{1}{10}$ . As in the previous section, we analyze the regret under the assumption that the good event holds. We start with the following regret decomposition,

$$\begin{aligned} R_K &= \sum_{k=1}^K \sum_{h;s} q_h(s) \left( \sum_{j=1}^D r_h^k(j; s) - \sum_{j=1}^D Q_h^k(j; s) \right) = \underbrace{\sum_{k=1}^K \sum_{h;s} q_h(s) \sum_{j=1}^D \left( r_h^k(j; s) - Q_h^k(j; s) \right)}_{\text{BIAS}_1} + \underbrace{\sum_{k=1}^K \sum_{h;s} q_h(s) \sum_{j=1}^D \left( Q_h^k(j; s) - Q_h^k(j; s) \right)}_{\text{BIAS}_2} + \underbrace{\sum_{k=1}^K \sum_{h;s} q_h(s) \sum_{j=1}^D \left( Q_h^k(j; s) - B_h^k(j; s) \right)}_{\text{BONUS}} \\ &+ \underbrace{\sum_{k=1}^K \sum_{h;s} q_h(s) \sum_{j=1}^D \left( B_h^k(j; s) - Q_h^k(j; s) \right)}_{\text{REG}} + \underbrace{\sum_{k=1}^K \sum_{h;s} q_h(s) \sum_{j=1}^D \left( Q_h^k(j; s) - B_h^k(j; s) \right)}_{\text{DRIFT}} \end{aligned}$$

where the first is by Lemma E.1.  $\text{BIAS}_2$  is bounded under event  $\mathcal{E}$  by  $O\left(\frac{H^2}{3}\right)$ , and  $\text{DRIFT}$  is bounded by  $O\left(H^5(K+D) + \frac{H^5 d_{\max}}{3}\right)$  by Lemma B.7. The other three terms are bounded in Lemmas C.4 to C.6. Overall,

$$\begin{aligned} R_K &\leq \frac{2}{3} \sum_{k=1}^K \sum_{h;s} q_h(s) b_h^k(s) + \sum_{k=1}^K \sum_{h;s} q_h(s) b_h^k(s) + O\left(H^4(K+D) + \frac{H^4 d_{\max}}{3} + \frac{H^2}{3}\right) \\ &+ O\left(\frac{H^2}{3}\right) + O\left(H^2 S A K + H^3 S^p \overline{AK} + H^4 S^3 A^2 + H^4 S^2 A d_{\max}\right) + \sum_{k=1}^K \sum_{h;s} q_h(s) b_h^k(s) \\ &+ \frac{H \ln A}{3} + \frac{1}{3} \sum_{k,h;s} q_h(s) b_h^k(s) + O\left(H^5 K + \frac{H^3}{2}\right) \\ &+ O\left(H^5(K+D) + \frac{H^5 d_{\max}}{3}\right) \\ &= O\left(\frac{H \ln A}{3} + H^2 S A K + H^5(K+D) + H^3 S^p \overline{AK} + \frac{H^3}{2} + \frac{H^2}{3} + \frac{H^5 d_{\max}}{3} + H^4 S^2 A d_{\max} + H^4 S^3 A^2\right) \end{aligned}$$

For  $\epsilon = \frac{1}{H^2 S A K + H^4 (K + D)}$  and  $\delta = 2H$  we get:

$$R_K \leq O\left(H^3 S^D \overline{AK} + H^3 D \overline{K + D} + H^4 S^2 A d_{\max} + H^4 S^3 A^2\right) \quad \square$$

### C.3. Bound on $\text{BIAS}_1$

Lemma C.4. Under the good event,

$$\text{BIAS}_1 \leq \frac{2}{3} \sum_{k=1}^K \sum_{h,s} q_h(s) b_h^k(s) + \sum_{k=1}^K \sum_{h,s} q_h(s) b_h^k(s) + O\left(H^4 (K + D) + \frac{H^4 d_{\max}}{3} + \frac{H^2}{3}\right) \quad \square$$

Proof. Let  $Y_k = \mathbb{P}_{h,s}^D \left[ q_h(s) \sum_{j=1}^{k+d^k} Q_h^k(j,s); Q_h^k(s) \right]$ . It holds that

$$\text{BIAS}_1 = \sum_{k=1}^K \sum_{h,s} q_h(s) \sum_{j=1}^{k+d^k} Q_h^k(j,s); Q_h^k(s) - \mathbb{E}_{k=1}^E [Y_k] + \sum_{k=1}^K \mathbb{E}_{k=1} [Y_k] - \sum_{k=1}^K Y_k$$

Under event  $\mathcal{E}^f$  it holds that  $\mathbb{P}_{k=1}^K \mathbb{E}_{k=1} [Y_k] = \mathbb{P}_{k=1}^K Y_k = \frac{1}{3} \mathbb{P}_{k=1}^K \sum_{h,s} q_h(s) b_h^k(s) + \frac{H^2}{3} \ln 10$ . In addition,

$$\begin{aligned} & \sum_{k=1}^K \sum_{h,s} q_h(s) \sum_{j=1}^{k+d^k} Q_h^k(j,s); Q_h^k(s) - \mathbb{E}_{k=1}^E [Y_k] = \\ &= \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) - \frac{q_h^k(s;a) r_h^k(s;a)}{q_h^k(s;a) +} \\ &= \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) - \frac{(q_h^k(s;a) + + q_h^k(s;a) - q_h^k(s;a)) r_h^k(s;a)}{q_h^k(s;a) +} \\ & \quad + \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) \frac{r_h^k(s;a)}{q_h^k(s;a) +} \\ & \quad + \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) \frac{r_h^k(s;a)(q_h^k(s;a) - q_h^k(s;a))}{q_h^k(s;a) +} \\ &= \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) \frac{(q_h^k(s;a) +)(1 - r_h^k(s;a))}{q_h^k(s;a) +} \\ & \quad \underbrace{\left\{ \frac{r_h^k(s;a)(q_h^k(s;a) - q_h^k(s;a))}{q_h^k(s;a) +} \right\}}_{(i)} \\ & \quad + \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) \frac{r_h^k(s;a)(q_h^k(s;a) - q_h^k(s;a))}{q_h^k(s;a) +} \\ & \quad \underbrace{\left\{ \frac{r_h^k(s;a)}{q_h^k(s;a) +} \right\}}_{(ii)} \\ & \quad + \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) \frac{r_h^k(s;a)}{q_h^k(s;a) +} \\ & \quad \underbrace{\left\{ \frac{r_h^k(s;a)(q_h^k(s;a) - q_h^k(s;a))}{q_h^k(s;a) +} \right\}}_{(iii)} \\ & \quad + \sum_{k=1}^K \sum_{h,s;a} q_h(s) \sum_{j=1}^{k+d^k} (a_j s) Q_h^k(s;a) \frac{r_h^k(s;a)(q_h^k(s;a) - q_h^k(s;a))}{q_h^k(s;a) +} \\ & \quad \underbrace{\left\{ \frac{r_h^k(s;a)(q_h^k(s;a) - q_h^k(s;a))}{q_h^k(s;a) +} \right\}}_{(iv)} \end{aligned}$$



where the first two inequalities are by definition  $q_h^k$  and the fact that  $\mathbb{P}^k$  under even  $E^p$ ; and the last inequality is since, by definition,  $q_h^k(s)$  maximize the probability to visit  $s$  at time  $h$  among all the occupancy measures with  $\mathbb{P}^k$ . Summing over  $k$  and bounding the first term above as follows completes the proof:

$$\begin{aligned} \sum_{k=1}^K \sum_{h;s} q_h^k(s) b_h^k(s) &= \sum_{k=1}^K \sum_{h;s;a} \frac{3H \sum_{h} q_h^k(s) \sum_{h}^{k+d^k} (a_j s) r_h^k(s; a)}{q_h^k(s; a) +} \\ &\quad + \sum_{k=1}^K \sum_{h;s;a} \frac{2H \sum_{h} q_h^k(s) \sum_{h}^{k+d^k} (a_j s) r_h^k(s; a) (q_h^k(s; a) - q_h^k(s; a))}{q_h^k(s; a) +} \\ &\leq 3H^2SAK + 2H \sum_{k=1}^K \sum_{h;s;a} (q_h^k(s; a) - q_h^k(s; a)) \\ &\leq O(H^2SAK + H^3S^pAK + H^4S^3A^2 + H^4S^2Ad_{\max}) ; \end{aligned}$$

where the last is by even  $E^p$ . □

### C.5. Bound on REG

Lemma C.6. For  $\frac{H \ln A}{2H}$  it holds that  $\text{REG} \leq \frac{H \ln A}{3} + \frac{1}{3} \sum_{k,h;s} q_h(s) b_h^k(s) + O(H^5K) + \frac{H^3}{2}$  :

Proof. By Corollary E.7, since  $\max_{k,h;s;a} B_h^k(s; a) \leq 5H^2$ ,

$$\begin{aligned} \text{REG} &\leq \frac{H \ln A}{2} + 2 \sum_{k,h;s;a} q_h(s) \sum_{h}^{k+d^k} (a_j s) \sum_{h} q_h^k(s; a) B_h^k(s; a)^2 + O(H^5K) \\ &\leq \frac{H \ln A}{2} + 2 \sum_{k,h;s;a} q_h(s) \sum_{h}^{k+d^k} (a_j s) \sum_{h} q_h^k(s; a)^2 + 2 \sum_{k,h;s;a} q_h(s) \sum_{h}^{k+d^k} (a_j s) B_h^k(s; a)^2 + O(H^5K) : (18) \end{aligned}$$

For the middle term

$$\begin{aligned} 2 \sum_{k,h;s;a} q_h(s) \sum_{h}^{k+d^k} (a_j s) \sum_{h} q_h^k(s; a)^2 &\leq 2 \sum_{k,h;s;a} q_h(s) \sum_{h}^{k+d^k} (a_j s) \frac{r_h^k(s; a)^2 H^2 \mathbb{1}_{s_h^k = s; a_h^k = a}}{(q_h^k(s; a) +)^2} \\ &\leq 2H^2 \sum_{k,h;s;a} q_h(s) \sum_{h}^{k+d^k} (a_j s) \frac{r_h^k(s; a)}{q_h^k(s; a) +} + O\left(\frac{H^3}{2}\right) \\ &\leq \frac{2}{3} H \sum_{k,h;s} q_h(s) b_h^k(s) + O\left(\frac{H^3}{2}\right) \\ &\leq \frac{1}{3} \sum_{k,h;s} q_h(s) b_h^k(s) + O\left(\frac{H^3}{2}\right) ; \end{aligned}$$

where the second inequality is  $E^p$  and the last is since  $\frac{H \ln A}{2H}$ . For the last term in Eq. (18) we use  $\sum_{k,h;s;a} B_h^k(s; a) \leq 5H$  and therefore  $\sum_{k,h;s;a} B_h^k(s; a)^2 \leq 5H^2$ . Thus,

$$\sum_{k,h;s;a} q_h(s) \sum_{h}^{k+d^k} (a_j s) B_h^k(s; a)^2 \leq 25H^5K :$$

Overall,

$$\text{REG} \leq \frac{H \ln A}{3} + \frac{1}{3} \sum_{k,h;s} q_h(s) b_h^k(s) + O(H^5K) + \frac{H^3}{2} : \quad \square$$



## D. Delay-Adapted Policy Optimization for Adversarial MDP with Linear Q-function

## Algorithm 5 Delay-Adapted Policy Optimization with Linear Q-function

Input: feature dimension  $m$ , action space  $A$ , horizon  $H$ , feature map  $f_h : S \times A \rightarrow \mathbb{R}^n$ , simulator of the environment, horizon  $H$ , learning rate  $\gamma > 0$ , exploration parameter  $\epsilon > 0$ , confidence parameter  $\delta > 0$ .

Initialization: Set approximation parameter  $\epsilon = (HnK)^{-1}$ , bonus parameters  $b_h^1 = H^{\frac{1}{2}} \frac{\delta}{n}$ ,  $b_h^2 = H^{\frac{1}{2}} \frac{\delta}{n}$ ,  $b_h^r = 2H^{\frac{1}{2}} \frac{\delta}{n}$ ,  $b_h^v = 4H^{\frac{1}{2}} \frac{\delta}{n}$ ,  $b_h^{if} = H^{\frac{1}{2}} \frac{\delta}{n}$  and Geometric Resampling parameters  $M = d^{\frac{24}{2}} \ln \frac{10H^2Kn}{\delta}$ ,  $N = d^2 \log \frac{1}{\delta}$ .

Define  $\mathbb{1}_h^j(a; s) = \mathbb{1}_{A^j}$  for every  $(s; a; h) \in S \times A \times [H]$ .

for  $k = 1; 2; \dots; K$  do

Play episode  $k$  with policy  $\pi^k$  and observe trajectory  $(s_h^k; a_h^k)_{h=1}^H$ .

# Policy Evaluation

for  $j$  such that  $j + d^j = k$  do

Observe bandit feedback  $(r_h^j; a_h^j)_{h=1}^H$ .

Collect  $MN$  trajectories  $T^j$  using the simulator and  $\mathbb{1}^j$ .

Compute estimated inverse covariance matrix  $\hat{\Sigma}_h^j = \text{GEOMETRICRESAMPLING}(T^j; M; N; \epsilon)$  using the Matrix Geometric Resampling procedure (Neu & Olkhovskaya, 2021) (which samples trajectories of the policy using the simulator).

Compute Monte-Carlo estimates  $\hat{c}_h^j = \frac{1}{M} \sum_{h^0=h}^H c_{h^0}^j(s_{h^0}^j; a_{h^0}^j)$  for every  $h \in [H]$ .

Compute estimated Q-function weights  $\hat{\Lambda}_h^j = \frac{1}{M} \sum_{h^0=h}^H \Lambda_{h^0}^j(s_{h^0}^j; a_{h^0}^j)$  for every  $h \in [H]$ .

Define delay-adapted ratio  $\hat{r}_h^j(s; a) = \frac{\hat{c}_h^j(s; a)}{\max_{a'} \hat{c}_h^j(s; a')}$  for every  $(s; a; h) \in S \times A \times [H]$ .

Define delay-adapted estimated Q-function  $\hat{Q}_h^j(s; a) = r_h^j(s; a) \hat{\Lambda}_h^j(s; a)$  for every  $(s; a; h) \in S \times A \times [H]$ .

Define bonus  $b_h^j(s; a) = b_h^{j:1}(s) + b_h^{j:2}(s; a) + b_h^{j:r}(s) + b_h^{j:v}(s; a) + b_h^{j:g}(s; a)$  for every  $(s; a; h)$ , where:

$$\begin{aligned} b_h^{j:1}(s) &= \frac{1}{\epsilon} \sum_{a'} r_h^j(s; a') \mathbb{1}_h^j(a; s) k_h(s; a) k_{\Lambda_{j:1}}^2(s; a); & b_h^{j:2}(s; a) &= \frac{2}{\epsilon} r_h^j(s; a) k_h(s; a) k_{\Lambda_{j:2}}^2(s; a); \\ b_h^{j:v}(s) &= \frac{1}{\epsilon} \sum_{a'} r_h^j(s; a') \mathbb{1}_h^j(a; s) k_h(s; a) k_{\Lambda_{j:v}}^2(s; a); & b_h^{j:r}(s; a) &= r_h^j(s; a) (1 - r_h^j(s; a)); \\ b_h^{j:f}(s) &= \frac{1}{\epsilon} \sum_{a'} r_h^j(s; a') \mathbb{1}_h^j(a; s) k_h(s; a) k_{\Lambda_{j:f}}^2(s; a); & b_h^{j:g}(s; a) &= g_h^j(s; a) k_h(s; a) k_{\Lambda_{j:g}}^2(s; a); \end{aligned}$$

end for

# Policy Improvement

Define the policy  $\pi^{k+1}$  for every  $(s; a; h) \in S \times A \times [H]$  by:

$$\pi_h^{k+1}(a; s) = \frac{\sum_{j: j+d^j=k} \exp\left(\frac{\hat{Q}_h^j(s; a) - \hat{B}_h^j(s; a)}{\epsilon}\right)}{\sum_{a' \in A} \sum_{j: j+d^j=k} \exp\left(\frac{\hat{Q}_h^j(s; a') - \hat{B}_h^j(s; a')}{\epsilon}\right)};$$

where  $\hat{B}_h^j(s; a)$  estimates  $B_h^j(s; a)$  using the bonus procedure in Algorithm 6.  
end for

Theorem D.1. Set  $\epsilon = \frac{\delta}{nK}$  and  $\gamma = \min\left\{\frac{\delta}{10Hd_{\max}}, \frac{1}{H(K+D)^{3+4}}\right\}$  and skip episodes with delay larger than  $\tau = D^{1+4}$ .

Running Algorithm 5 in an adversarial MDP  $\mathcal{M} = (S; A; H; p; f; c^k; g_{k=1}^k)$  with linear Q-function, access to a simulator of the environment, a known features map  $f_h : S \times A \rightarrow \mathbb{R}^n$  and delays  $d^k; g_{k=1}^k$  guarantees, with probability at least  $1 - \delta$ ,

$$R_K \leq O\left(H^3 n^{5+4} K^{3+4} \log \frac{KHA}{\delta} + H^2 D^{3+4} \log \frac{KHA}{\delta} + H^5 n \log \frac{KHA}{\delta}\right);$$

Remark D.2. As noted in the main text, Algorithm 6 is not sample efficient, and may require  $(\frac{KHA}{\delta})^{O(H)}$  calls to the simulator. However, under the stronger assumption that the MDP is linear (e.g., see assumption 2.1 in Sherman et al. (2023)), we can replace this procedure with the OLSPE procedure of Sherman et al. (2023) which would make our algorithm fully efficient while achieving the same regret guarantee of Theorem D.1.

Algorithm 6 Bonus Procedure (Algorithm 3 in Luo et al. (2021) adapted to non-dilated bonuses)

Input: episode  $j$ , horizon  $h$ , states, actions, local bonus function  $b_h$ , simulator of the environment.  
 if  $B_h^j(s; a)$  was computed before then  
     Return the previously computed value  $B_h^j(s; a)$ .  
 end if  
 Compute  $j_h(j; s)$  and  $j_h^{j+d^j}(j; s)$  (which involves recursive calls to Bonus Procedure to compute  $B_h^{j^0}(j^0; s^0) < j + d^j$ ).  
 Sample  $s^0 \sim p_h(j; s; a)$  using the simulator.  
 Compute  $j_{h+1}(j; s^0)$  and sample  $a^0 \sim j_{h+1}(j; s^0)$ .  
 Return  $j_h(j; a) + B_{h+1}^j(s^0; a^0)$ .

D.1. The good event

Let  $\epsilon = 10 \log \frac{10KH}{\delta}$ ,  $H^k$  be the history of all episodes  $\{j : j + d^j < k\}$  and define  $E_k[\cdot] = E[\cdot | H^k]$ . Let  $\| \cdot \|_{op}$  be the operator norm. That is, give a matrix  $A \in \mathbb{R}^{n \times n}$ ,  $\|A\|_{op} := \inf_{\|x\|_2=1} \|Ax\|_2$ . Define the following events:

$$E^8 = \left\{ \begin{aligned} & \|E_k\|_{op} \leq \frac{1}{2} \text{ and } \|E_k\|_{op} \leq \frac{1}{2} \text{ for } k \leq K; \\ & \|E_k\|_{op} \leq \frac{1}{2} \text{ for } k \leq K; \end{aligned} \right.$$

$$E^b = \left\{ \begin{aligned} & \sum_{k=1}^K \sum_a E_{s, q_h}^{k+d^k}(a; j; s) m^{k+d^k} Q_h^k(s; a)^2 \leq \frac{4H^2 d_{\max} \log \frac{10H}{\delta}}{2}; \end{aligned} \right.$$

$$E^f = \left\{ \begin{aligned} & \sum_{k=1}^K \sum_h E_{s, q_h}^{k+d^k}(j; s; Q_h^k(s; )) \leq \frac{H^2 \log \frac{10}{\delta}}{2} + O(H^2 K); \end{aligned} \right.$$

$$E^g = \left\{ \begin{aligned} & \sum_{k=1}^K \sum_h E_{s, q_h}^{k+d^k}(j; s; Q_h^k(s; )) \leq \frac{H^2 \log \frac{10}{\delta}}{2} + O(H^2 K); \end{aligned} \right.$$

$$E^B = \left\{ \begin{aligned} & \sum_{k=1}^K \sum_{h=1}^K E_{s, q_h}^{k+d^k}(j; s; B_h^k(s; )) \leq O(H^3 \frac{1}{nK} \ln \frac{10}{\delta}); \end{aligned} \right.$$

The good event is the intersection of the above events. The following lemma establishes that the good event holds with high probability.

Lemma D.3 (The Good Event) Let  $G = E^8 \cap E^b \cap E^f \cap E^g \cap E^B$  be the good event. It holds that  $\mathbb{P}[G] \geq 1 - \delta$ .

Proof. We'll show that each of the events  $E^8; E^b; E^f; E^g; E^B$  occurs with probability  $\geq 1 - \delta/5$ . By taking the

union bound we'll get that  $\Pr[G] \leq 10^{-1}$ .

Event E<sup>a</sup>: We use GEOMETRICRESAMPLING with  $M = \frac{24}{\epsilon} \ln \frac{10H^2Kn}{\epsilon}$  and  $N = \frac{2}{\epsilon} \ln \frac{1}{\epsilon}$ . Thus, by Lemma D.5, the event holds for each  $k$  and  $h$  separately with probability of at least  $1 - \frac{\epsilon}{10HK}$ . By taking the union bound over  $[H]$  and  $[K]$  we get that  $\Pr[E^a] \leq 10^{-1}$ .

Event E<sup>b</sup>: Fix  $h$ . By Lemma D.4,  $|\hat{Q}_h^k(s; a) - Q_h^k(s; a)| \leq \frac{H}{\sqrt{m^{k+d^k}}}$  and thus  $\sum_{a \in \mathcal{A}} E_{s \sim q_h} \sum_{j \in \mathcal{S}} m^{k+d^k} (a; j, s) \sum_{a \in \mathcal{A}} \hat{Q}_h^k(s; a) \leq \frac{H^2 d_{\max}^2}{2}$ .

Also note that  $m^{k+d^k}$  is determined by the history  $\mathcal{H}^{k+d^k}$ . Thus, by Lemma E.4 the event holds with probability  $1 - \frac{\epsilon}{10H}$ . The proof is finished by a union bound over  $[H]$ .

Event E<sup>f</sup>: Let  $Y_k = \sum_{h=1}^D E_{s \sim q_h} \sum_{j \in \mathcal{S}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a)$ . Similarly to the tabular case, we use a variant of Freedman's inequality (Lemma E.3) to bound  $\sum_{k=1}^K X_k := \sum_{k=1}^K E_k[Y_k] - \sum_{k=1}^K Y_k$ .

$$\begin{aligned} E_k X_k^2 &= E_k Y_k^2 = E_k \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (a; j, s) \hat{Q}_h^k(s; a) \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \\ &= E_k \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \\ &\leq H E_k \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \\ &\leq H E_k \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \\ &= H E_k \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \\ &= H \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) E_k \hat{Q}_h^k(s; a) \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) \hat{Q}_h^k(s; a) \\ &\leq H^3 \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) r_h^k(s; a) k_h(s; a) k_{h^k}^2 + O(H^4) \end{aligned}$$

where the first inequality is since  $\|x\|_1^2 \leq n \|x\|_2^2$  for any  $x \in \mathbb{R}^n$ , the second inequality is by Jensen's inequality, and the last inequality is by Lemma D.10. Also  $\sum_{k=1}^K Y_k \leq \frac{H^2}{\epsilon}$ . Therefore by Lemma E.3 with probability  $1 - \frac{\epsilon}{10}$ ,

$$\begin{aligned} \sum_{k=1}^K E_k[Y_k] - \sum_{k=1}^K Y_k &\leq H \sum_{k=1}^K \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) r_h^k(s; a) k_h(s; a) k_{h^k}^2 + \frac{H^2 \log 10}{\epsilon} + O(HK) \\ &= \sum_{k=1}^K \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) r_h^k(s; a) k_h(s; a) k_{h^k}^2 + \frac{H^2 \log 10}{\epsilon} + O(H^2 K) \end{aligned}$$

Event E<sup>g</sup>: The proof is similar to event E<sup>f</sup>.

Event E<sup>B</sup>: Define,

$$X_k := \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) r_h^k(s; a) k_h(s; a) k_{h^k}^2 - \sum_{h=1}^D \sum_{j \in \mathcal{S}} \sum_{a \in \mathcal{A}} m^{k+d^k} (j; s) r_h^k(s; a) k_h(s; a) k_{h^k}^2$$

and note that  $\sum_{k=1}^K X_k \leq O(H^3 D \bar{n})$ . By Azuma–Hoeffding inequality (Lemma E.2) we get that the event holds with probability  $1 - \frac{\epsilon}{10}$ .  $\square$

Lemma D.4. For any  $(k; h; s; a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , it holds that  $\|\hat{Q}_h^k(s; a)\| \leq \frac{H}{\sqrt{2}}$ .

Proof. By Cauchy-Schwartz, Eq. (32) in Lemma E.11 and Assumption 2.1:

$$\|\hat{Q}_h^k(s; a)\| = \|\mathbb{E}_{s \sim q_h} [r_h^k(s; a) + \gamma \mathbb{E}_{s' \sim q_h} [V_h^k(s') - V_h^k(s)] | s, a]\| \leq H \sqrt{\mathbb{E}_{s \sim q_h} [k_2^2 + \mathbb{E}_{s' \sim q_h} [k_2^2]]} \leq H \sqrt{2} = \frac{H}{\sqrt{2}}.$$

□

Lemma D.5 (Lemma D.1 in Luo et al. (2021)) Fix a policy  $\pi$  with a covariance matrix  $\Sigma_h = \mathbb{E}_{s \sim q_h} [ (s; a) (s; a)^\top ]$ . Let  $\hat{\Sigma}_h^+$  be the output of GEOMETRICRESAMPLING( $T; M; N; \cdot$ ) with  $M = \frac{24}{\epsilon} \ln \frac{nH}{\delta}$  and  $N = \frac{2}{\epsilon} \ln \frac{1}{\delta}$  and  $T$  are  $MN$  trajectories collected with  $\pi$ . Then,  $\|\hat{\Sigma}_h^+ - \Sigma_h\|_{\text{op}} \leq \epsilon$ , and with probability  $1 - \delta$ ,

$$\|\hat{\Sigma}_h^+ - \Sigma_h\|_{\text{op}} \leq \epsilon; \quad \|\hat{\Sigma}_h^+ - \Sigma_h\|_{\text{op}} \leq \epsilon + 2\epsilon;$$

## D.2. Proof of the main theorem

Proof of Theorem D.1 By Lemma D.3, the good event holds with probability of at least  $1 - \delta$ . As in the previous section, we analyze the regret under the assumption that the good event holds. We start with the following regret decomposition,

$$\begin{aligned} R_K &= \sum_{k=1}^K \mathbb{E}_{s \sim q_h} [r_h^k(j; s) - \mathbb{E}_{s \sim q_h} [r_h^k(j; s) | s]] + \sum_{k=1}^K \mathbb{E}_{s \sim q_h} [r_h^k(j; s) - \mathbb{E}_{s \sim q_h} [r_h^k(j; s) | s]] \\ &= \underbrace{\sum_{k=1}^K \mathbb{E}_{s \sim q_h} [r_h^k(j; s) - \mathbb{E}_{s \sim q_h} [r_h^k(j; s) | s]]}_{\text{BIAS}_1} + \underbrace{\sum_{k=1}^K \mathbb{E}_{s \sim q_h} [r_h^k(j; s) - \mathbb{E}_{s \sim q_h} [r_h^k(j; s) | s]]}_{\text{BIAS}_2} \\ &\quad + \underbrace{\sum_{k=1}^K \mathbb{E}_{s \sim q_h} [r_h^k(j; s) - \mathbb{E}_{s \sim q_h} [r_h^k(j; s) | s]]}_{\text{REG}} + \underbrace{\sum_{k=1}^K \mathbb{E}_{s \sim q_h} [r_h^k(j; s) - \mathbb{E}_{s \sim q_h} [r_h^k(j; s) | s]]}_{\text{BONUS}} \\ &\quad + \underbrace{\sum_{k=1}^K \mathbb{E}_{s \sim q_h} [r_h^k(j; s) - \mathbb{E}_{s \sim q_h} [r_h^k(j; s) | s]]}_{\text{DRIFT}}; \end{aligned}$$

where the first is by Lemma E.1. The five terms above are bounded in Lemmas D.6, D.8, D.9, D.11 and D.13. For  $\frac{1}{10Hd_{\max}}$  and  $\frac{1}{(HnK)^4}$ ,

$$\begin{aligned}
 R_K & \leq \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s, q_h} \left[ b_h^{k;1}(s) + b_h^{k;f}(s) + O\left(-H^3 \bar{n}(K+D)^p + \frac{H^2}{\bar{n}} + H^p \bar{n}K + H^2 K\right) \right]}_{\text{BIAS}_1} \\
 & + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s; a, q_h} \left[ b_h^{k;2}(s; a) + b_h^{k;r}(s; a) + b_h^{k;g}(s; a) + O\left(\frac{H^2}{\bar{n}} + H^p \bar{n}K + H^2 K\right) \right]}_{\text{BIAS}_2} \\
 & + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s, q_h} \left[ b_h^{k;v}(s) + O\left(\frac{H}{\bar{n}} + H^5 nK + \frac{H^2}{\bar{n}} + H^3\right) \right]}_{\text{REG}} \quad \underbrace{\left| \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s, q_h} \left[ b_h^k(s; a) \right] \right|}_{\text{DRIFT}} \\
 & + O\left(\frac{H^p}{\bar{n}} H^2 nK + H^2 nK + -H^3 \bar{n}(K+D) + H^3 \bar{n}K\right) + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s; a, q_h} \left[ b_h^k(s; a) \right]}_{\text{BONUS}} \\
 & \leq O\left(\frac{H}{\bar{n}} + H^5 nK + -H^4 \bar{n}(K+D) + \frac{H^p}{\bar{n}} H^2 nK + H^3 \bar{n}K\right) :
 \end{aligned}$$

For  $\frac{1}{10Hd_{\max}}$  and  $\frac{1}{(HnK)^4}$ , we get:

$$R_K \leq O\left(H^3 n^{5=4} K^{3=4} + H^2 D^{3=4} + H^4 nK^{1=4} + H^3 \bar{n}K + H^2 d_{\max} \frac{H^p}{\bar{n}K} \right) :$$

Finally, by skipping rounds larger than  $\frac{1}{10Hd_{\max}}$ , we get:  $R_K \leq O\left(H^3 n^{5=4} K^{3=4} + H^2 D^{3=4} + H^5 n\right)$   $\square$

### D.3. Bound on $\text{BIAS}_1$

Lemma D.6. Under the good event,

$$\text{BIAS}_1 \leq \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s, q_h} \left[ b_h^{k;1}(s) + b_h^{k;f}(s) + O\left(-H^3 \bar{n}(K+D)^p + \frac{H^2}{\bar{n}} + H^p \bar{n}K + H^2 K\right) \right] :$$

Proof. We first decompose  $\text{BIAS}_1$  as,

$$\begin{aligned}
 \text{BIAS}_1 & = \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s, q_h} \left[ b_h^{k+d^k}(j; s); Q_h^k(s; \cdot) - \hat{Q}_h^k(s; \cdot) \right] \\
 & = \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s, q_h} \left[ b_h^{k+d^k}(j; s); E_k \left[ Q_h^k(s; \cdot) - \hat{Q}_h^k(s; \cdot) \right] \right]}_{\text{(A)}} \\
 & + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s, q_h} \left[ b_h^{k+d^k}(j; s); E_k \left[ \hat{Q}_h^k(s; \cdot) - Q_h^k(s; \cdot) \right] \right]}_{\text{(B)}} ;
 \end{aligned}$$

where the (B)  $\sum_{k=1}^P \sum_{h;a} \mathbb{E}_{s \sim q_h} [b_h^{k,f}(s)] + \frac{H^2}{k} + O(H^2 K)$  by event  $\mathcal{E}^f$ . For (A) we use Lemma D.7,

$$\begin{aligned}
 (A) \quad & \sum_{k=1}^P \sum_{h;a} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i \\
 & + \sum_{k=1}^P \sum_{h;a} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (1 - r_h^k(s; a)) (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i + O(H^2 K) \\
 = & \underbrace{\sum_{k=1}^P \sum_{h;a} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i}_{(i)} \\
 & + \underbrace{\sum_{k=1}^P \sum_{h;a} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (1 - r_h^k(s; a)) (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i}_{(ii)} \\
 & + \underbrace{\sum_{k=1}^P \sum_{h;a} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (1 - r_h^k(s; a)) (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i}_{(iii)} + O(H^2 K);
 \end{aligned}$$

Using Cauchy-Schwarz inequality,

$$\begin{aligned}
 (i) \quad & \sum_{h;a,k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i \\
 & \leq H^P \frac{1}{n} \sum_{h;a,k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i \\
 & \leq H^P \frac{1}{n} \sum_{h;a,k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i + O(HK^P \frac{1}{n}) \\
 = & \sum_{h;k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i + O(HK^P \frac{1}{n});
 \end{aligned}$$

where the second inequality is since  $\frac{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} = \frac{q}{k_h^k (1 + \frac{k}{h})^{-1} \frac{k}{h}} \frac{r}{k_h^k k^2} H^q \frac{1}{n}$  by Eq. (31) of Lemma E.11 and Assumption 2.1, and the third is by Cauchy-Schwarz and event  $\mathcal{E}^f$  in a similar way to Eq. (21). For term (ii), note that

$$\frac{h^{k+d^k} (a_j s)_h (1 - r_h^k(s; a)) (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \leq \frac{1}{k_h^k k^2} \frac{H^P \frac{1}{n}}{k_h^k k^2}; \quad (19)$$

where the first inequality is by Cauchy-Schwarz, and the second inequality is by Eq. (32) in Lemma E.11. Therefore,

$$\begin{aligned}
 (ii) \quad & H^P \frac{1}{n} \sum_{h;a,k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h (1 - r_h^k(s; a)) (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i \\
 = & H^P \frac{1}{n} \sum_{h;a,k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h \frac{\max_{j \in \mathcal{A}} \{h^{k+d^k} (a_j s)_h; \frac{k}{h} (a_j s)_h\}}{\max_{j \in \mathcal{A}} \{h^{k+d^k} (a_j s)_h; \frac{k}{h} (a_j s)_h\}}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i \\
 & \leq H^P \frac{1}{n} \sum_{h;a,k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h \frac{\max_{j \in \mathcal{A}} \{h^{k+d^k} (a_j s)_h; \frac{k}{h} (a_j s)_h\}}{\max_{j \in \mathcal{A}} \{h^{k+d^k} (a_j s)_h; \frac{k}{h} (a_j s)_h\}}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i \\
 & \leq H^P \frac{1}{n} \sum_{h;k} \mathbb{E}_{s \sim q_h} \left[ \frac{h^{k+d^k} (a_j s)_h \frac{\max_{j \in \mathcal{A}} \{h^{k+d^k} (a_j s)_h; \frac{k}{h} (a_j s)_h\}}{\max_{j \in \mathcal{A}} \{h^{k+d^k} (a_j s)_h; \frac{k}{h} (a_j s)_h\}}}{h^{k+d^k} (a_j s)_h (s; a)^> (1 + \frac{k}{h})^{-1} \frac{k}{h}} \right]^i \\
 & \leq O\left(-H^3 \frac{1}{n} (K + D)\right);
 \end{aligned}$$

where the last is by Lemma D.12. Similarly, using Eq. (33) in Lemma E.11,

$$j_h(s; a) > (I + \frac{k}{h})^{-1} \frac{k}{h} k_{hj} k_h(s; a) k_2 k (I + \frac{k}{h})^{-1} \frac{k}{h} k_2 k_h k_2 H^P \bar{n}. \quad (20)$$

Thus, (iii) is bounded in the same way as (ii).  $\square$

Lemma D.7. Under the good event, for every  $(k; h; s; a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , it holds that

$$E_k Q_h^k(s; a) - \hat{Q}_h^k(s; a) = j_h(s; a) > (I + \frac{k}{h})^{-1} \frac{k}{h} + (1 - r_h^k(s; a)) j_h(s; a) > (I + \frac{k}{h})^{-1} \frac{k}{h} k_h k_h O(H):$$

Proof. By the definition of  $Q_h^k(s; a)$ ,  $\hat{Q}_h^k(s; a)$  and  $\hat{\Lambda}_h^k$ ,

$$\begin{aligned} E_k Q_h^k(s; a) - \hat{Q}_h^k(s; a) &= j_h(s; a) > \frac{k}{h} r_h^k(s; a) E_k \hat{\Lambda}_h^{k;+} \\ &= j_h(s; a) > \frac{k}{h} r_h^k(s; a) E_k \hat{\Lambda}_h^{k;+} - E_k j_h(s_h^k; a_h^k) L_h^j \\ &= j_h(s; a) > \frac{k}{h} r_h^k(s; a) (I + \frac{k}{h})^{-1} E_k j_h(s_h^k; a_h^k) L_h^j + O(H) \\ &= j_h(s; a) > \frac{k}{h} r_h^k(s; a) (I + \frac{k}{h})^{-1} E_k j_h(s_h^k; a_h^k) E L_h^j j_h^k; a_h^k + O(H) \\ &= j_h(s; a) > \frac{k}{h} r_h^k(s; a) (I + \frac{k}{h})^{-1} E_k j_h(s_h^k; a_h^k) Q_h^j(s_h^k; a_h^k) + O(H) \\ &= j_h(s; a) > \frac{k}{h} r_h^k(s; a) (I + \frac{k}{h})^{-1} E_k j_h(s_h^k; a_h^k) j_h(s_h^k; a_h^k) > \frac{k}{h} + O(H) \\ &= j_h(s; a) > \frac{k}{h} r_h^k(s; a) (I + \frac{k}{h})^{-1} \frac{k}{h} k_h k_h + O(H) \\ &= j_h(s; a) > (I + \frac{k}{h})^{-1} (I + \frac{k}{h})^{-1} \frac{k}{h} r_h^k(s; a) (I + \frac{k}{h})^{-1} \frac{k}{h} k_h k_h + O(H) \\ &= j_h(s; a) > (I + \frac{k}{h})^{-1} \frac{k}{h} + (1 - r_h^k(s; a)) j_h(s; a) > (I + \frac{k}{h})^{-1} \frac{k}{h} k_h k_h + O(H) \end{aligned}$$

where the third equality is since,

$$\begin{aligned} j_h(s; a) > r_h^k(s; a) (E_k \hat{\Lambda}_h^{k;+} - (I + \frac{k}{h})^{-1} E_k j_h(s_h^k; a_h^k) L_h^j) \\ &= k_h(s; a) k_2 (E_k \hat{\Lambda}_h^{k;+} - (I + \frac{k}{h})^{-1} E_k j_h(s_h^k; a_h^k) L_h^j) \\ &= H (\hat{\Lambda}_h^{k;+} - (I + \frac{k}{h})^{-1} E_k j_h(s_h^k; a_h^k) L_h^j) \leq 2H; \end{aligned} \quad (21)$$

where the first inequality is by Cauchy–Schwarz, the second inequality is by Assumption 2.1 and the last is by exact Assumption 2.1. In the same way we have,

$$j_h(s; a) > (I + \frac{k}{h})^{-1} \frac{k}{h} + (1 - r_h^k(s; a)) j_h(s; a) > (I + \frac{k}{h})^{-1} \frac{k}{h} k_h k_h E_k Q_h^k(s; a) - \hat{Q}_h^k(s; a) = O(H): \quad \square$$

#### D.4. Bound on $\text{BIAS}_2$

Lemma D.8. Under the good event,

$$\text{BIAS}_2 = \sum_{k=1}^K \sum_{h=1}^H E_{s; a} j_h(s; a) (b_h^{k;2}(s; a) + b_h^{k;r}(s; a) + b_h^{k;g}(s; a)) + O\left(\frac{H^2}{n} + H^P \bar{n} K + H^2 K\right);$$

Proof. Similarly to  $\text{BIAS}_1$ , we first decompose,

$$\begin{aligned} \text{BIAS}_2 &= \sum_{k=1}^K \sum_{h=1}^H E_{s; a} j_h(s; a) (Q_h^k(s; a) - \hat{Q}_h^k(s; a)) \\ &= \underbrace{\sum_{k=1}^K \sum_{h=1}^H E_{s; a} j_h(s; a) (E_k Q_h^k(s; a) - \hat{Q}_h^k(s; a))}_{(A)} + \underbrace{\sum_{k=1}^K \sum_{h=1}^H E_{s; a} j_h(s; a) (E_k \hat{Q}_h^k(s; a) - \hat{Q}_h^k(s; a))}_{(B)}; \end{aligned}$$

where (B)  $\mathbb{P}_{k,h} \mathbb{E}_{s;a} [b_h^{k,g}(s;a)] + \frac{H^2}{k} + O(H^2 K)$  by event  $\mathcal{E}^g$ , and for (A) we use again Lemma D.7,

$$\begin{aligned}
 (A) \quad & \sum_{k=1}^K \sum_{h;a} \mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} \\
 & + \sum_{k=1}^K \sum_{h;a} \mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} + O(H^2 K) \\
 = & \sum_{k=1}^K \sum_{h;a} \underbrace{\mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h}}_{(i)} \\
 & + \sum_{k=1}^K \sum_{h;a} \underbrace{\mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h}}_{(ii)} \\
 & + \sum_{k=1}^K \sum_{h;a} \underbrace{\mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} + O(H^2 K)}_{(iii)}
 \end{aligned}$$

Similar to (i) in Lemma D.6

$$\begin{aligned}
 (i) \quad & \sum_{h;a} \mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} \\
 & \leq \frac{1}{H^p \bar{n}} \sum_{h;a} \mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} \\
 & \leq \frac{1}{H^p \bar{n}} \sum_{h;a} \mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} + O(H^p \bar{n}) \\
 = & \sum_{h;a} \mathbb{E}_{s;a} [b_h^{k,2}(s;a)] + O(H^p \bar{n}):
 \end{aligned}$$

Using Eq. (19) and Eq. (20),

$$\begin{aligned}
 (ii) \quad & \frac{1}{H^p \bar{n}} \sum_{h;a} \mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} = \frac{1}{2} \sum_{h;a} \mathbb{E}_{s;a} [b_h^{k,r}(s;a)] \\
 (iii) \quad & \frac{1}{H^p \bar{n}} \sum_{h;a} \mathbb{E}_{s; q_h} [r_h^k(s;a) - \hat{r}_h^k(s;a)]^2 (1 + \frac{k}{h})^{-1} \frac{k}{h} = \frac{1}{2} \sum_{h;a} \mathbb{E}_{s;a} [b_h^{k,r}(s;a)] : \quad \square
 \end{aligned}$$

### D.5. Bound on REG

Lemma D.9. Under the good event, for  $\frac{1}{10Hd_{\max}}$ ,

$$\text{REG} \leq \sum_{k=1}^K \sum_h \mathbb{E}_{s; q_h} [b_h^{k,v}(s)] + O(H^5 n K) + \frac{H \log A}{H} + \frac{H^2}{H} + H^3 :$$



Proof. Note that  $(\hat{Q}_h^k(s; a) - \hat{B}_h^k(s; a)) = \hat{Q}_h^k(s; a) - \hat{B}_h^k(s; a)$ . Thus, using lemma Lemma E.8 for each  $(h, s)$ ,

$$\begin{aligned} & \sum_{k=1}^K \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} (\hat{Q}_h^k(s; a) - \hat{B}_h^k(s; a))^2 \\ & \leq \frac{\log A}{2} + 2 \sum_{k=1}^K \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} \hat{Q}_h^k(s; a)^2 + 2 \sum_{k=1}^K \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} \hat{B}_h^k(s; a)^2. \end{aligned} \quad (22)$$

Note that  $\hat{B}_h^k(s; a) = O(H^D \bar{n})$ . Thus  $\hat{B}_h^k(s; a) = O(H^{2D} \bar{n})$ , and the last sum can be bounded by,

$$\sum_{k;a} m^{k+d^k} \hat{B}_h^k(s; a)^2 = O(H^{4n}) \sum_{k;a} m^{k+d^k} = O(H^{4n}) \sum_{k=1}^K m^{k+d^k} = O(H^4 n K). \quad (23)$$

Taking the expectation with respect to  $q_h$  on the first sum in (22), by event  $\mathbb{E}^b$ ,

$$\begin{aligned} & 2 \sum_{k=1}^K \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} (\hat{Q}_h^k(s; a) - \hat{B}_h^k(s; a))^2 \\ & \leq 4 \sum_{k=1}^K \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} E_k \hat{Q}_h^k(s; a)^2 + \frac{8 H^2 d_{\max}}{2} \\ & \leq 4 \sum_{k=1}^K \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} E_k \hat{Q}_h^k(s; a)^2 + \frac{H}{2}; \end{aligned} \quad (24)$$

where the last is since  $\frac{1}{10 H d_{\max}}$ . Finally, by Lemma D.10,

$$\begin{aligned} & 4 \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} E_k \hat{Q}_h^k(s; a)^2 \leq 4 H^2 \sum_a \sum_{j \in \mathcal{S}} m^{k+d^k} r_h^k(s; a) k_h(s; a) k_h^2(s; a) + O(H^2); \\ & = b_h^{k;v}(s) + O(H^2); \end{aligned} \quad (25)$$

Combining the above with Eqs. (22) to (25), summing over  $k$  and taking the expectation completes the proof.  $\square$

Lemma D.10. Under event  $\mathbb{E}^b$ , for every  $(k; h; s; a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , it holds that

$$E_k \hat{Q}_h^k(s; a)^2 \leq H^2 r_h^k(s; a) k_h(s; a) k_h^2(s; a) + O(H^2);$$

Proof. By the definition of  $\hat{Q}_h^k(s; a)$  and  $\hat{B}_h^k(s; a)$ ,

$$\begin{aligned} E_k \hat{Q}_h^k(s; a)^2 & = E_k \left( r_h^k(s; a) - \sum_{i \in \mathcal{S}} \frac{h(s; a) \hat{Q}_h^k(s; a)}{h(s; a)} \right)^2 \\ & \leq H^2 E_k \left( r_h^k(s; a) - \sum_{i \in \mathcal{S}} \frac{h(s; a) \hat{Q}_h^k(s; a)}{h(s; a)} \right)^2 \\ & = H^2 \left( r_h^k(s; a) - \sum_{i \in \mathcal{S}} \frac{h(s; a) \hat{Q}_h^k(s; a)}{h(s; a)} \right)^2 \\ & = H^2 \left( r_h^k(s; a) - \sum_{i \in \mathcal{S}} \frac{h(s; a) \hat{Q}_h^k(s; a)}{h(s; a)} \right)^2 \\ & = H^2 \left( r_h^k(s; a) - \sum_{i \in \mathcal{S}} \frac{h(s; a) \hat{Q}_h^k(s; a)}{h(s; a)} \right)^2 \end{aligned} \quad (26)$$

We rewrite,

$$\begin{aligned}
 h(s; a)^{\wedge k;+} k_h^{\wedge k;+} h(s; a) &= \underbrace{h(s; a)^{\wedge k;+} k_h^{\wedge k;+} (1 + \frac{k}{h})^{-1} h(s; a)}_{(i)} \\
 &+ \underbrace{h(s; a)^{\wedge k;+} (1 + \frac{k}{h})^{-1} k_h (1 + \frac{k}{h})^{-1} h(s; a)}_{(ii)} \\
 &+ \underbrace{h(s; a)^{\wedge k;+} (1 + \frac{k}{h})^{-1} k_h^k (1 + \frac{k}{h})^{-1} h(s; a)}_{(iii)} : \tag{27}
 \end{aligned}$$

We now bound each of the above as follows:

$$\begin{aligned}
 (i) \quad & h(s; a)^{\wedge k;+} k_h^{\wedge k;+} (1 + \frac{k}{h})^{-1} h(s; a) \quad \text{(Cauchy-Schwarz)} \\
 & k_h h(s; a) k_2^{\wedge k;+} k_h^{\wedge k;+} (1 + \frac{k}{h})^{-1} k_h h(s; a) k_2 \\
 & \frac{\wedge k;+}{h} k_h^{\wedge k;+} (1 + \frac{k}{h})^{-1} \quad \text{(} k_h h(s; a) k_2^{-1} \text{)} \\
 & (1 + \frac{k}{h})^2 ; \tag{28}
 \end{aligned}$$

where the last inequality is by event  $\mathcal{E}_t$ . Similarly,

$$\begin{aligned}
 (ii) \quad & k_h h(s; a) k_2^{\wedge k;+} (1 + \frac{k}{h})^{-1} k_h (1 + \frac{k}{h})^{-1} h(s; a) \quad \text{(2)} \\
 & \frac{\wedge k;+}{h} (1 + \frac{k}{h})^{-1} \quad \text{(} k_h (1 + \frac{k}{h})^{-1} h(s; a) \text{)} \\
 & 2 \frac{k}{h} (1 + \frac{k}{h})^{-1} h(s; a) \quad \text{(2)} \\
 & 2 k_h h(s; a) k_2 \quad \text{(2)} ; \tag{29}
 \end{aligned}$$

where the third inequality is by event  $\mathcal{E}_t$  and the fourth inequality uses Eq. (33) in Lemma E.11. Finally, using Eq. (35) in Lemma E.11,

$$\begin{aligned}
 (iii) \quad & h(s; a)^{\wedge k;+} (1 + \frac{k}{h})^{-1} h(s; a) \\
 & = h(s; a)^{\wedge k;+} (1 + \frac{k}{h})^{-1} \wedge h(s; a) + k_h h(s; a) k_h^{\wedge k;+} \\
 & k_h h(s; a) k_2 (1 + \frac{k}{h})^{-1} \wedge h(s; a) \quad \text{(2)} + k_h h(s; a) k_h^{\wedge k;+} \\
 & (1 + \frac{k}{h})^{-1} \wedge h(s; a) \quad \text{(2)} + k_h h(s; a) k_h^{\wedge k;+} \\
 & 2 + k_h h(s; a) k_h^{\wedge k;+} : \tag{30}
 \end{aligned}$$

Combining Eqs. (26) to (30) completes the proof.  $\square$

#### D.6. Bound on DRIFT

Lemma D.11. If  $\frac{1}{H} \leq \frac{1}{\bar{n}}$  and  $\frac{P}{4Hd_{\max}}$  then,  $\text{DRIFT} \leq -H^4 \bar{n} (K + D)$  :

Proof. Note that if  $\frac{1}{H} \leq \frac{1}{\bar{n}}$  then  $k_h(s; a) \leq (H^P \bar{n})$  and thus  $B_h^k(s; a) \leq (H^{2P} \bar{n})$ . Now, using Lemma D.12,

$$\begin{aligned}
 \text{DRIFT} &= \sum_{k=1}^K \sum_h \mathbb{E}_{s \sim q_h} \left[ \frac{hD}{k_h(j; s)} \frac{k_h^{k+d^k}(j; s)}{k_h^{k+d^k}(j; s)} Q_h^k(s; \cdot) \mathbb{E}_i B_h^k(s; \cdot) \right] \\
 &\leq (H^{2P} \bar{n}) \sum_{k=1}^K \sum_h \mathbb{E}_{s \sim q_h} \left[ k_h^k(j; s) \frac{k_h^{k+d^k}(j; s)}{k_h^{k+d^k}(j; s)} k_1^i \right] \leq -H^4 \bar{n} (K + D) : \quad \square
 \end{aligned}$$

Lemma D.12. If  $\frac{1}{H^2 \bar{n}}$  and  $\frac{P}{4Hd_{\max}}$  then for each  $(h; s)$ :  $\sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}[\mathbb{B}_h^k(s; a)] = O(-H(K+D))$ :

Proof. We apply lemma Lemma E.9 for each  $(k; h; s; a)$  with  $\hat{a} = \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{Q}_h^j(s; a)$  and  $M = \frac{10H}{4Hd_{\max}}$  since  $\mathbb{Q}_h^k(s; a) \leq \frac{H}{6H^2 \bar{n}}$  by Lemma D.4, and  $\mathbb{B}_h^k(s; a) \leq \frac{6H}{4Hd_{\max}}$  whenever  $\frac{1}{H^2 \bar{n}}$  which implies that  $\mathbb{B}_h^k(s; a) \leq \frac{6H}{4Hd_{\max}}$ . We get that,

$$\begin{aligned} \sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}[\mathbb{B}_h^k(s; a)] &= \sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}[\mathbb{B}_h^k(s; a)] \\ &= \sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}[\mathbb{Q}_h^j(s; a^0) + \mathbb{B}_h^j(s; a^0)] + \frac{10H}{4Hd_{\max}} \\ &+ \sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}[\mathbb{Q}_h^j(s; a) + \mathbb{B}_h^j(s; a)] + \frac{10H}{4Hd_{\max}} \\ &= \sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}[\mathbb{Q}_h^j(s; a^0) + \mathbb{B}_h^j(s; a^0)] + \frac{20H}{4Hd_{\max}} + \sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}[\mathbb{Q}_h^j(s; a) + \mathbb{B}_h^j(s; a)] + \frac{20H}{4Hd_{\max}} \\ &= \sum_{k=1}^K \sum_{j: k+d^k \leq j} \sum_{i: j+d^i \leq k+d^k} \frac{40H}{4Hd_{\max}}; \end{aligned}$$

Summing the above over  $k$  and applying Lemma E.10 completes the proof. □

### D.7. Bound on BONUS

Lemma D.13. Under the good event,

$$\text{BONUS} = \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)] + O\left(\frac{P}{H^2} 2nK + H^2 nK + -H^3 \frac{P}{\bar{n}}(K+D) + H^3 \frac{P}{nK}\right);$$

Proof. Under event  $\mathcal{E}^B$ ,

$$\begin{aligned} \sum_{h; k} \mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)] &= \sum_{h; k} \mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)] + O\left(H^3 \frac{P}{nK}\right) \\ &= \sum_{h; k} \mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)] + O\left(H^3 \frac{P}{nK}\right); \end{aligned}$$

where the equality is by Lemma D.14. Recall that  $b_h^k(s; a) = b_h^{k;1}(s) + b_h^{k;2}(s; a) + b_h^{k;v}(s) + b_h^{k;r}(s; a) + b_h^{k;f}(s) + b_h^{k;g}(s; a)$ . Now, the expectation  $\mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)]$  of each of the bonus terms is bounded in Lemmas D.15 to D.20. □

Lemma D.14. For any  $k$ , let  $b_h^k(s; a)$  be a loss function determined by the history  $\mathcal{H}_h^{k+d^k}$ , and let  $\mathbb{B}_h^k(s; a)$  be a randomized bonus function such that, for every  $(h; s; a) \in \mathcal{H}_h^{k+d^k}$ ,

$$\mathbb{E}_k[\mathbb{B}_h^k(s; a)] = b_h^k(s; a) + \mathbb{E}_{s^0} \mathbb{P}_h(j; s; a) \mathbb{E}_{a^0} \sum_{i: j+d^i \leq k+d^k} \mathbb{E}_k[\mathbb{B}_{h+1}^k(s^0; a^0)]$$

$$\text{Then, } \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)] = \sum_{h=1}^H \mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)] + \sum_{h=1}^H \mathbb{E}_{s; a} \mathbb{E}_{q_h} [b_h^k(s; a)];$$



Proof. For any  $h$  and  $k$ ,

$$E_{s;a} \sum_{q_h^k} b_h^{k;2}(s; a) = \frac{1}{2} E_{s;a} \sum_{q_h^k} r_h^k(s; a) k_h(s; a) k_h^{\wedge k;+} + \frac{1}{2} E_{s;a} \sum_{q_h^k} k_h^k(s; a) k_h^{\wedge k;+} \frac{1}{2} \sqrt{\frac{1}{n(1+2)}};$$

where the last is as in the proof of Lemma D.15. □

Lemma D.17. Under the good event,  $\prod_{k=1}^K \prod_{h=1}^H E_{s;a} \sum_{q_h^k} b_h^{k;r}(s; a) = O(-H^3 D \sqrt{n(K+D)})$  :

Proof. For any  $h$  and  $k$ ,

$$\begin{aligned} E_{s;a} \sum_{q_h^k} b_h^{k;r}(s; a) &= \frac{1}{r} E_{s;a} \sum_{q_h^k} r_h^k(s; a) = \frac{1}{r} E_s \sum_{q_h^k} \sum_a X_{a}^k(j; s) (1 - r_h^k(s; a)) \\ &= \frac{1}{r} E_s \sum_{q_h^k} \sum_a X_{a}^k(j; s) \frac{\max_{a'} k_h^k(a; s); k_h^{k+d^k}(a; s) g_{a'}^k(j; s)}{\max_{a'} k_h^k(a; s); k_h^{k+d^k}(a; s) g_{a'}^k(j; s)} \\ &= \frac{1}{r} E_s \sum_{q_h^k} \sum_a \max_{a'} k_h^k(a; s); k_h^{k+d^k}(a; s) g_{a'}^k(j; s) \\ &= \frac{1}{r} E_s \sum_{q_h^k} k_h^{k+d^k}(j; s) k_h^k(j; s) k_1^i : \end{aligned}$$

Finally, taking the sum and applying Lemma D.12 completes the proof. □

Lemma D.18. Under the good event,  $\prod_{k=1}^K \prod_{h=1}^H E_{s;a} \sum_{q_h^k} b_h^{k;v}(s) = O((1+2)^H n K)$  :

Proof. For any  $h$  and  $k$ ,

$$\begin{aligned} E_s \sum_{q_h^k} b_h^{k;v}(s) &= \frac{1}{v} m^{k+d^k} E_s \sum_{q_h^k} \sum_a X_{a}^{k+d^k}(j; s) r_h^k(s; a) k_h(s; a) k_h^{\wedge k;+} \\ &= \frac{1}{v} m^{k+d^k} E_s \sum_{q_h^k} \sum_a X_{a}^k(j; s) k_h(s; a) k_h^{\wedge k;+} \\ &= \frac{1}{v} m^{k+d^k} E_{s;a} \sum_{q_h^k} k_h(s; a) k_h^{\wedge k;+} \frac{1}{v} m^{k+d^k} n(1+2) : \end{aligned}$$

Summing over  $h$  and  $k$  and noting that  $\prod_{k=1}^K m^{k+d^k} = K$  completes the proof. □

Lemma D.19. Under the good event,  $\prod_{k=1}^K \prod_{h=1}^H E_{s;a} \sum_{q_h^k} b_h^{k;f}(s) = O((1+2)^H n K)$  :

Proof. Similarly to Lemma D.18  $E_s \sum_{q_h^k} b_h^{k;f}(s) = \frac{1}{f} n(1+2)$ . □

Lemma D.20. Under the good event,  $\prod_{k=1}^K \prod_{h=1}^H E_{s;a} \sum_{q_h^k} b_h^{k;g}(s; a) = O((1+2)^H n K)$  :

Proof. Again, similarly to Lemma D.18  $E_{s;a} \sum_{q_h^k} b_h^{k;g}(s; a) = \frac{1}{g} n(1+2)$  : □

### E. Auxiliary Lemmas

Lemma E.1 (Value Difference Lemma (Even-Dar et al., 2009)) For any two policies  $\pi_1$  and  $\pi_2$ ,

$$V_{\pi_1}(s_{\text{init}}) - V_{\pi_2}(s_{\text{init}}) = \sum_{h=1}^H \mathbb{E}_{s \sim q_h} [v_h(\pi_1(s)) - v_h(\pi_2(s); Q_h(s; \pi_1))]$$

Lemma E.2 (Azuma–Hoeffding inequality) Let  $\{X_t\}_{t=1}^T$  be a real valued martingale difference sequence adapted to a filtration  $F_1 \subseteq F_2 \subseteq \dots$  (i.e.,  $\mathbb{E}[X_t | F_t] = 0$ ). If  $|X_t| \leq R$  a.s. then with probability at least  $1 - \delta$ ,

$$\sum_{t=1}^T X_t \leq R \sqrt{T \ln \frac{1}{\delta}}$$

Lemma E.3 (A special form of Freedman's Inequality, Theorem 1 of Beygelzimer et al. (2011)) Let  $\{X_t\}_{t=1}^T$  be a real valued martingale difference sequence adapted to a filtration  $F_1 \subseteq F_2 \subseteq \dots$  (i.e.,  $\mathbb{E}[X_t | F_t] = 0$ ). If  $|X_t| \leq R$  a.s. then for any  $\epsilon \in (0, 1-R)$ ;  $T \geq 2N/\epsilon^2$  it holds with probability at least  $1 - \delta$ ,

$$\sum_{t=1}^T X_t \leq \sum_{t=1}^T \mathbb{E}[X_t^2 | F_t] + \frac{\log(1/\delta)}{\epsilon}$$

Lemma E.4 (Consequence of Freedman's Inequality, e.g., Lemma E.2 in (Cohen et al., 2021)) Let  $\{X_t\}_{t=1}^T$  be a sequence of random variables, supported  $[-R, R]$ , and adapted to a filtration  $F_1 \subseteq F_2 \subseteq \dots$ . For any  $T$ , with probability  $1 - \delta$ ,

$$\sum_{t=1}^T X_t \leq 2\mathbb{E}[X_t | F_t] + 4R \sqrt{T \ln \frac{1}{\delta}}$$

Lemma E.5 (Lemma A.2 of Luo et al. (2021)) Given a filtration  $F_0 \subseteq F_1 \subseteq \dots$ , let  $z_h^k(s; a) \in [0, R]$  and  $c_h^k(s; a) \in [0, 1]$  be sequences of  $\mathcal{F}_k$ -measurable functions.  $\mathbb{I}_h^k(s; a) \in [0, R]$  is a sequence of random variables such that  $\mathbb{E}[\mathbb{I}_h^k(s; a) | F_k] = z_h^k(s; a)$  then with probability  $1 - \delta$ ,

$$\sum_{k=1}^K \sum_{h; s; a} \frac{\mathbb{I}_h^k(s; a) - c_h^k(s; a) z_h^k(s; a)}{c_h^k(s; a)} \leq \sum_{k=1}^K \sum_{h; s; a} \frac{c_h^k(s; a) z_h^k(s; a)}{c_h^k(s; a)} \leq \frac{RH}{2} \ln \frac{H}{\delta}$$

Lemma E.6 (Lemma 9 of Thune et al. (2019)) Let  $\epsilon > 0$ , variables delays  $d_k \in [0, K]$ , and loss vectors  $\mathbf{a}^k \in [0, 1]^A$  for all  $k \in [K]$ . Define,

$$\mathbf{1}^k(\mathbf{a}) = \frac{1}{A}; \quad \mathbf{1}^{k+1}(\mathbf{a}) = \frac{\prod_{j: j+d_j=k} \mathbf{a}^j}{\prod_{j: j+d_j=k} \mathbf{a}^j(\mathbf{a}^0)}$$

Then, for any  $\delta \in (0, 1/A)$ :

$$\sum_{k=1}^K \sum_{\mathbf{a} \in \mathcal{A}} (\mathbf{1}^{k+d_k}(\mathbf{a}) - \mathbf{1}^k(\mathbf{a})) \leq \frac{\ln A}{\delta} + \sum_{k=1}^K \sum_{\mathbf{a} \in \mathcal{A}} \mathbf{1}^{k+d_k}(\mathbf{a}) \mathbf{1}^k(\mathbf{a})^2$$

Corollary E.7. Let  $\epsilon > 0$ , variables delays  $d_k \in [0, K]$ , and loss vectors  $\mathbf{a}^k \in [0, 1]^A$  for all  $k \in [K]$ . Define,

$$\mathbf{1}^k(\mathbf{a}) = \frac{1}{A}; \quad \mathbf{1}^{k+1}(\mathbf{a}) = \frac{\prod_{j: j+d_j=k} \mathbf{a}^j}{\prod_{j: j+d_j=k} \mathbf{a}^j(\mathbf{a}^0)}$$

Then, for any  $\delta \in (0, 1/A)$ :

$$\sum_{k=1}^K \sum_{\mathbf{a} \in \mathcal{A}} (\mathbf{1}^{k+d_k}(\mathbf{a}) - \mathbf{1}^k(\mathbf{a})) \leq \frac{\ln A}{\delta} + 2 \sum_{k=1}^K \sum_{\mathbf{a} \in \mathcal{A}} \mathbf{1}^{k+d_k}(\mathbf{a}) \mathbf{1}^k(\mathbf{a})^2 + 2KM^2$$

Proof. Note that,

$$w^1(a) = \frac{1}{A} ; \quad w^{k+1}(a) = \frac{e^{-\sum_{j:j+d^j=k} \ell^j(a) + M}}{e^{-\sum_{j:j+d^j=k} \ell^j(a^0) + M}}$$

The statement now follows immediately by applying Lemma E.6 on the losses  $\ell^k$ . □

Lemma E.8. Let  $\delta > 0$ , variables delays  $d_k^k$ , and loss vectors  $\ell^k \in [0, 1]^A$  for all  $k \in [K]$ . Define,

$$w^k(a) = \frac{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a)\right)}{e^{-\sum_{j:j+d^j < k} \ell^j(a^0)}};$$

where the empty sum is zero. If  $\sum_{j:j+d^j=k} \ell^k(a) > \delta$  for all  $k \in [K]$ , then,

$$\sum_{k=1}^K \sum_a w^{k+d^k}(a) \ell^k(a) \leq \frac{\ln A}{\delta} + \sum_{k=1}^K \sum_a w^{k+d^k}(a) m^{k+d^k} \delta^2;$$

where  $m^k = \sum_{j:j+d^j=k} \ell^k(a)$ .

Proof. The proof is based in part on the proof of Thune et al. (2019, Lemma 9). Define  $w^k(a) = \frac{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a)\right)}{e^{-\sum_{j:j+d^j < k} \ell^j(a^0)}}$ . We have that

$$\begin{aligned} \frac{w^{k+1}}{w^k} &= \frac{\exp\left(-\sum_{j:j+d^j < k+1} \ell^j(a)\right)}{e^{-\sum_{j:j+d^j < k+1} \ell^j(a^0)}} \frac{e^{-\sum_{j:j+d^j < k} \ell^j(a^0)}}{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a)\right)} \\ &= \frac{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a) - \ell^k(a)\right)}{e^{-\sum_{j:j+d^j < k} \ell^j(a^0) - \ell^k(a^0)}} \frac{e^{-\sum_{j:j+d^j < k} \ell^j(a^0)}}{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a)\right)} \\ &= \frac{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a) - \ell^k(a)\right)}{e^{-\sum_{j:j+d^j < k} \ell^j(a^0) - \ell^k(a^0)}} \frac{e^{-\sum_{j:j+d^j < k} \ell^j(a^0)}}{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a)\right)} \\ &= \frac{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a) - \ell^k(a)\right)}{e^{-\sum_{j:j+d^j < k} \ell^j(a^0) - \ell^k(a^0)}} \frac{e^{-\sum_{j:j+d^j < k} \ell^j(a^0)}}{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a)\right)} \\ &= \exp\left(-\sum_{j:j+d^j=k} \ell^k(a) + \sum_{j:j+d^j=k} \ell^k(a^0)\right) \frac{e^{-\sum_{j:j+d^j < k} \ell^j(a^0)}}{\exp\left(-\sum_{j:j+d^j < k} \ell^j(a)\right)} \end{aligned}$$

where the first inequality is since  $1 + x + x^2 \leq e^x$  for  $x \leq 1$ , the second inequality is since  $e^{-x} \leq 1 - nx + \frac{x^2}{2}$  for any  $x \in \mathbb{R}^n$ , and the last inequality is since  $e^{-x} \leq e^x$ . Telescoping the ratio above for  $k=1, \dots, K$  we get,

$$\frac{w^{K+1}}{w^1} \leq \exp\left(-\sum_{k=1}^K \sum_a w^{k+d^k}(a) \ell^k(a) + \sum_{k=1}^K \sum_a w^{k+d^k}(a) m^{k+d^k} \delta^2\right);$$

where we used that for all  $j, j + d^j \leq K$  (we can assume w.l.o.g that all the missing feedback is observed in the end of the interaction). On the other hand,

$$\begin{aligned}
 \frac{W^{K+1}}{W^1} &\geq \frac{\sum_a \exp\left(-\eta \sum_{j:j+d^j \leq K} \ell^j(a)\right)}{A} \\
 &\geq \frac{\max_a \exp\left(-\eta \sum_{j:j+d^j \leq K} \ell^j(a)\right)}{A} \\
 &\geq \frac{\exp\left(-\min_a \eta \sum_{j:j+d^j \leq K} \ell^j(a)\right)}{A} \\
 &\geq \frac{\exp\left(-\eta \sum_{j:j+d^j \leq K} \sum_a \pi(a) \ell^j(a)\right)}{A} \\
 &= \frac{\exp\left(-\eta \sum_{k=1}^K \sum_a \pi(a) \ell^k(a)\right)}{A},
 \end{aligned}$$

where again we used that for all  $j, j + d^j \leq K$  for the last equality. Combining the last two inequalities taking  $\ln$  on both sides and rearranging the terms we get,

$$\sum_{k=1}^K \sum_a \left(\pi^{k+d^k}(a) - \pi(a)\right) \ell^k(a) \leq \frac{\ln A}{\eta} + \eta \sum_{k=1}^K \sum_a \pi^{k+d^k}(a) m^{k+d^k} (\ell^k(a))^2. \quad \square$$

**Lemma E.9** (Lemma 1 of Cesa-Bianchi et al. (2019) adapted to negative losses). *Let  $\pi, \tilde{\pi} \in \Delta_A$  and  $\ell \in [-M, \infty)^A$  such that*

$$\tilde{\pi}(a) = \frac{\pi(a) e^{-\eta \ell(a)}}{\sum_{a^\theta} \pi(a^\theta) e^{-\eta \ell(a^\theta)}}.$$

*It holds that,*

$$-\eta \pi(a) (\ell(a) + M) \leq \tilde{\pi}(a) - \pi(a) \leq \eta \tilde{\pi}(a) \sum_{a^\theta} \pi(a^\theta) (\ell(a^\theta) + M).$$

*Proof.* By the condition in the lemma,

$$\begin{aligned}
 \tilde{\pi}(a) &= \frac{\pi(a) \exp(-\eta \ell(a))}{\sum_{a^\theta} \pi(a^\theta) \exp(-\eta \ell(a^\theta))} \\
 &= \frac{\pi(a) \exp(-\eta(\ell(a) + M))}{\sum_{a^\theta} \pi(a^\theta) \exp(-\eta(\ell(a^\theta) + M))} \\
 &\geq \frac{\pi(a) \exp(-\eta(\ell(a) + M))}{\sum_{a^\theta} \pi(a^\theta)} \\
 &= \pi(a) \exp(-\eta(\ell(a) + M)) \\
 &\geq \pi(a) (1 - \eta(\ell(a) + M)),
 \end{aligned}$$

where in the first inequality we use the fact that  $\ell(a^\theta) + M \geq 0$  and so the exponent at the denominator  $\leq 1$ ; and the second inequality is by  $e^{-x} \geq 1 - x$ . Thus,

$$\tilde{\pi}(a) - \pi(a) \geq -\eta \pi(a) (\ell(a) + M).$$

Similarly,

$$\begin{aligned}
 \tilde{\pi}(a) &= \frac{\pi(a) \exp(-\eta \ell(a))}{\sum_{a^\theta} \pi(a^\theta) \exp(-\eta \ell(a^\theta))} \\
 &= \frac{\pi(a) \exp(-\eta(\ell(a) + M))}{\sum_{a^\theta} \pi(a^\theta) \exp(-\eta(\ell(a^\theta) + M))} \\
 &\leq \frac{\pi(a)}{\sum_{a^\theta} \pi(a^\theta) \exp(-\eta(\ell(a^\theta) + M))},
 \end{aligned}$$



where the inequality is since  $\ell(a) + M \geq 0$ . Thus,

$$\begin{aligned} \tilde{\pi}(a) - \pi(a) &\leq \tilde{\pi}(a) \left( 1 - \sum_{a^\theta} \pi(a^\theta) \exp(-\eta(\ell(a^\theta) + M)) \right) \\ &= \tilde{\pi}(a) \sum_{a^\theta} \pi(a^\theta) (1 - \exp(-\eta(\ell(a^\theta) + M))) \\ &\leq \eta \tilde{\pi}(a) \sum_{a^\theta} \pi(a^\theta) (\ell(a^\theta) + M). \end{aligned} \quad \square$$

**Lemma E.10** (Thune et al. (2019)). *Let  $\{d^k\}_{k=1}^K$  be a sequence of non-negative delays such that  $\sum_{k=1}^K d^k = D$ . Then,*

$$\sum_{k=1}^K \sum_{i=1}^K |\{k \leq i + d^i < k + d^k\}| \leq D + K.$$

*Proof.* The proof appears as part of the proof of Theorem 1 in Thune et al. (2019) or as a separate lemma in Jin et al. (2022, Lemma C.7)  $\square$

**Lemma E.11.** *If  $A$  is a positive semi-definite (PSD) matrix and  $\gamma > 0$ , then for any vector  $x \in \mathbb{R}^d$ ,*

$$x^T (A + \gamma I)^{-1} x \leq \frac{1}{\gamma} \|x\|_2^2 \quad (31)$$

$$x^T (A + \gamma I)^{-2} x \leq \frac{1}{\gamma^2} \|x\|_2^2 \quad (32)$$

$$x^T (A + \gamma I)^{-1} A x \leq \|x\|_2^2 \quad (33)$$

$$x^T (A + \gamma I)^{-1} A^2 (A + \gamma I)^{-1} x \leq \|x\|_2^2 \quad (34)$$

$$x^T (A + \gamma I)^{-1} A (A + \gamma I)^{-1} x \leq x^T (A + \gamma I)^{-1} x \quad (35)$$

*Proof.* By the spectral decomposition  $A = U^T D U$  where  $U$  is orthogonal matrix and  $D = \text{diag}(\lambda_1, \dots, \lambda_d)$  with  $\lambda_i \geq 0$ . Note that,

$$A + \gamma I = U^T (D + \gamma I) U \implies (A + \gamma I)^{-1} = U^T (D + \gamma I)^{-1} U.$$

Denote  $y = Ux$ . Then,

$$x^T (A + \gamma I)^{-1} x = y^T (D + \gamma I)^{-1} y = \sum_{i=1}^d \frac{1}{\lambda_i + \gamma} y_i^2 \leq \sum_{i=1}^d \frac{1}{\gamma} y_i^2 = \frac{1}{\gamma} \|y\|_2^2 = \frac{1}{\gamma} \|x\|_2^2,$$

which establishes Eq. (31). Eq. (32) is done similarly by noting that  $(A + \gamma I)^{-2} = U^T (D + \gamma I)^{-2} U$  since  $U^T U = I$ . For Eq. (33),

$$(A + \gamma I)^{-1} A = U^T (D + \gamma I)^{-1} D U.$$

Hence,

$$x^T (A + \gamma I)^{-1} A x = y^T (D + \gamma I)^{-1} D y = \sum_{i=1}^d \frac{\lambda_i}{\lambda_i + \gamma} y_i^2 \leq \sum_{i=1}^d y_i^2 = \|y\|_2^2 = \|x\|_2^2.$$

Eq. (34) is done similarly. For Eq. (35),

$$(A + \gamma I)^{-1} A (A + \gamma I)^{-1} = U^T (D + \gamma I)^{-1} D (D + \gamma I)^{-1} U.$$

Thus,

$$\begin{aligned}
 x^T(A + \gamma I)^{-1}A(A + \gamma I)^{-1}x &= y^T(D + \gamma I)^{-1}D(D + \gamma I)^{-1}y = \sum_{i=1}^d \frac{\lambda_i}{(\lambda_i + \gamma)^2} y_i^2 \\
 &\leq \sum_{i=1}^d \frac{1}{\lambda_i + \gamma} y_i^2 = y^T(D + \gamma I)^{-1}y \\
 &= x^T U^T (D + \gamma I)^{-1} U x = x^T (A + \gamma I)^{-1} x. \quad \square
 \end{aligned}$$

## F. DAPPO Implementation Details and Additional Experiments

### F.1. DAPPO Implementation Details

Our experiments are based on the implementation of PPO from the Stable-Baselines3 library (Raffin et al., 2021). All of the implementation details remain identical to the original implementation (including the architecture of the Deep Neural Networks and the default hyper-parameters), except for the two following modifications: (i) The objective of DAPPO, and (ii) We mimic learning with delayed feedback by withholding feedback from the algorithm for  $d$  steps.

PPO maintains a policy network  $\pi^\theta$  and a value network  $V^\phi$ . In each round  $k$ , the algorithm collects a rollout  $(s_1^k, a_1^k, r_1^k, \dots, s_H^k, a_H^k, r_H^k)$  of length  $H = 2048$  (notice that for the experiments we switched from costs to rewards). Note that the rollout is of length  $H$ , regardless of the number of episodes it takes to fill the rollout buffer to be of that length. That is, if for example the environment has a termination state and the episode ends before time  $H$ , we start a new episode and keep filling the buffer until we reach  $H$  interactions of the policy with the environment. That way, we can emulate fixed finite horizon MDPs as in our setting, even if the environment is not of fixed horizon. To this end, we treat each rollout of length  $H$  as a single episode. Apart from the state, action, reward and next state, the rollout buffer also stores the probability to take the chosen action  $\pi^{\theta^k}(a_h^k | s_h^k)$  for each  $h \in [H]$ .

Now, since we want to simulate delayed feedback, we do not use the buffer of round  $k$  to update  $\pi^{\theta^k}$ . Instead, we store this buffer and load the buffer from round  $k - d$ , where  $d$  is the delay in terms of episodes and not in terms of timesteps. I.e., if the delay in terms of timesteps is  $\tilde{d}$ , then  $d = \lfloor \tilde{d}/H \rfloor$ . At this point, the policy network objectives for DPPO and DAPPO are  $L_D^k(\theta)$  and  $L_{DA}^k(\theta)$ , respectively, where,

$$L_D^k(\theta) = \sum_{h=1}^H \min \left\{ g_h^k(\theta) \hat{A}_h^{k-d}, \text{clip}_{1-\epsilon} \left( g_h^k(\theta) \right) \hat{A}_h^{k-d} \right\};$$

$$L_{DA}^k(\theta) = \sum_{h=1}^H \min \left\{ R_h^k(\theta) \hat{A}_h^{k-d}, \text{clip}_{1-\epsilon} \left( R_h^k(\theta) \right) \hat{A}_h^{k-d} \right\},$$

for  $g_h^k(\theta) = \frac{\pi^k(a_h^k | s_h^k)}{\pi^k(a_h^k | s_h^k)}$  and  $R_h^k(\theta) = \frac{\pi^k(a_h^k | s_h^k)}{\max_{\pi^k} \pi^k(a_h^k | s_h^k)}$ .  $\hat{A}_h^{k-d} = L_h^{k-d} - V^\phi(s_h^{k-d})$  is an estimate of the advantage function, where  $L_h^{k-d}$  is the realized cost-to-go from  $(s_h^{k-d}, a_h^{k-d})$  until the first termination state in the rollout buffer. Note that  $L_D^k(\theta)$  is computed solely based on the parameters  $\theta$  and on data stored in rollout buffer, and does not require any modification to the the original algorithm. On the other hand, DAPPO computes, in addition, the probabilities of the last policy  $\pi^{\theta^k}$  over the trajectory of  $\pi^{\theta^{k-d}}$  (which has a relatively small computational cost).

The value network is trained simply by optimizing the mean-squared error (MSE) loss  $\sum_{h=1}^H (L_h^{k-d} - V^\phi(s_h^{k-d}))^2$ . Finally, the optimization of both the policy network and the value network is done using the Adam optimizer (Kingma & Ba, 2014) with learning rate  $\eta = 0.0003$ , batch size of 64 and for 10 epochs over the rollout buffer (these parameters remain unchanged from the original implementation of Raffin et al. (2021)).

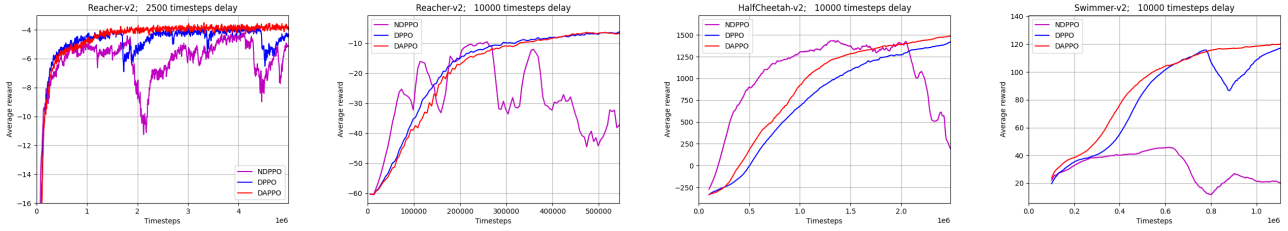


Figure 3. **Instability of NDPPO:** The learning curve of NDPPO vs DAPPO and DPPO under various settings. Plots show reward over a single run. x-axis is the number of timesteps until massive drop in performance (or up to 5M)

## F.2. Additional Experiments – Instability of NDPPO

We conducted experiments to show that NDPPO is an unstable algorithm. Note that  $\pi^{\theta^{k-d}}(a_h^{k-d} | s_h^{k-d})$  is not likely to be very small since  $a_h^{k-d}$  was sampled from  $\pi^{\theta^{k-d}}(\cdot | s_h^{k-d})$ . On the other hand,  $\pi^{\theta^k}(a_h^k | s_h^k)$  may be effectively 0 if the probability to choose  $a_h^k$  has decreased dramatically between time  $k-d$  and time  $k$ . This emphasizes that NDPPO (described in Section 5) does not only optimize over a biased objective, but is also highly unstable since the gradient of  $L_{ND}^k(\theta)$  is inversely proportional to  $\pi^{\theta^k}(a_h^k | s_h^k)$ . To demonstrate this phenomena, we present in Fig. 3 a few runs of NDPPO (compared to DPPO and DAPPO) under various settings. In some cases, such as in the SWIMMER-v2 environment, NDPPO is not able to improve the policy due to the large bias of its estimator. In other cases, such as the REACHER-v2 or HALFCHEETAH-v2 environments, the learning curve initially behaves similarly to DAPPO and DPPO, and sometimes even slightly better. This is due to the fact that  $\pi^{\theta^k}(a_h^k | s_h^k)$  is likely to be smaller than  $\pi^{\theta^{k-d}}(a_h^{k-d} | s_h^{k-d})$ , leading to larger updates (compared to DPPO and DAPPO). However, at some point, the NDPPO’s learning curve becomes much noisier due to the dramatic updates whenever  $\pi^{\theta^k}(a_h^k | s_h^k) \approx 0$ . This may result in an unrecoverable drop in performance even when the delay is small (as presented in Fig. 3), and in general, it gives an unstable algorithm with huge variance. DAPPO naturally avoids this issue by taking the maximum between  $\pi^{\theta^k}(a_h^k | s_h^k)$  and  $\pi^{\theta^{k-d}}(a_h^{k-d} | s_h^{k-d})$ .

## F.3. Additional Experiments – Drop in Performance as Delay Length Increases

We conducted experiments to exhibit the drop in performance that occurs when delay length increases. Fig. 4 compares the training curves of DPPO vs DAPPO with different lengths of  $\tilde{d} \in \{10000, 25000, 50000, 75000, 100000\}$ , alongside the training curve of PPO without delay. As expected, when the delay is relatively small (e.g.,  $\tilde{d} = 10000$ ), there is no significant difference between learning with or without delayed feedback. As the delay becomes larger, the performance of all algorithms drops (but at different rates).

One exception is Hal fCheetah-v2 where even for small delay performance drops. However, this is mainly due to the fact that the performance across different seeds is very noisy. Specifically, in 3 out of the 5 seeds, PPO without delays converges to a local maxima which has average reward of  $\approx 1500$  (similar to DPPO and DAPPO with sufficiently large delay). In the two other seeds it converges to a much better policy, which explains the large std in these graphs.

Whenever the delay becomes sufficiently large ( $\sim 25K - 50K$ ), there is a massive drop in the performance. This emphasize the great challenge that online algorithms need to face in the presence of delays. Namely, the algorithm updates its current policy based on estimated advantage function of a very different policy than the current one. This is also the point where the way we estimate the advantage function becomes important and the difference between DPPO and DAPPO becomes much more significant.

