KL-REGULARIZATION IS SUFFICIENT IN CONTEXTUAL BANDITS AND RLHF

Anonymous authors

000

001

002003004

006

008 009

010 011

012

013

014

016

018

019

021

024

025

026

027

028

029

031

033

037

038

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Recently, reinforcement learning from human feedback (RLHF) has demonstrated remarkable efficiency in fine-tuning large language models (LLMs), fueling a surge of interest in KL regularization. Yet, the theoretical foundations of KL regularization remain underexplored. Many prior works employ either explicit online exploration strategies—such as UCB, Thompson sampling, and forced sampling—or optimism-embedded optimization techniques (e.g., Xie et al. 2024) in addition to KL regularization to achieve sublinear regret in online RLHF. In this paper, we show, for the first time to our best knowledge, that such additional exploration strategies are unnecessary if KL regularization is already included. That is, KL regularization alone suffices to guarantee sublinear regret. We propose KL-EXP (and its RLHF variant, OEPO), an algorithm that achieves logarithmic KL-regularized regret—the standard objective in KL-regularized contextual bandits and RLHF—while also attaining $\tilde{\mathcal{O}}(\sqrt{T})$ unregularized regret, both under general function approximation. As a special case, in linear contextual bandits, we establish a $\mathcal{O}(\sqrt{dT \log N})$ bound on the unregularized regret, where d is the feature dimension and N is the number of arms. To our best knowledge, this is the first $\mathcal{O}(\sqrt{dT \log N})$ -type regret bound achieved without resorting to supLin-type algorithms, making it substantially more practical. Our experiments on linear and neural bandits, as well as on LLM fine-tuning with RLHF, demonstrate that our algorithms significantly outperform the baselines while remaining practical.

1 Introduction

The Kullback–Leibler (KL)-regularized contextual bandit problem (Langford & Zhang, 2007; Neu et al., 2017; Xiong et al., 2023; Xie et al., 2024) has recently attracted considerable attention due to its remarkable empirical success in fine-tuning large language models (LLMs), an application commonly referred to as reinforcement learning from human feedback (RLHF) (Christiano et al., 2017; Bai et al., 2022; Ouyang et al., 2022). This framework uses KL regularization as a key mechanism to balance reward optimization with distributional stability.

Despite these practical successes, the theoretical understanding of KL-regularization remains limited, particularly in the context of online learning. *Online exploration* is crucial for efficiently gathering informative feedback and addressing user preferences in RLHF. In this vein, many prior works have leveraged additional mechanisms to promote exploration, such as Upper Confidence Bound (UCB) (Xiong et al., 2023; 2024; Zhao et al., 2025), forced sampling (Zhao et al., 2024), and value-incentivized policy optimization (Xie et al., 2024; Cen et al., 2024). Building on these strategies, Xiong et al. (2023); Ye et al. (2024); Xie et al. (2024); Xiong et al. (2024); Cen et al. (2024) established $\mathcal{O}(\sqrt{T})$ bounds on *KL-regularized regret* (or $\mathcal{O}(1/\epsilon^2)$ sample complexity). More recently, Zhao et al. (2024; 2025) achieved the first logarithmic KL-regularized regret (or $\mathcal{O}(1/\epsilon)$ sample complexity).

However, optimizing the KL-regularized objective (Equation 1) already yields a randomized policy of the Gibbs distribution form (Equation 2). This implies that KL-regularization induces inherent exploration. Therefore, a natural question arises:

Can <u>logarithmic KL-regularized regret be achieved without extra exploration techniques</u> in contextual bandits and RLHF when KL-regularization is used?

Beyond this, we raise a more fundamental question: is achieving sublinear *KL-regularized regret*, by itself, truly sufficient? To the best of our knowledge, the tightest bound to date is $\mathcal{O}(\eta \log T)$, established by Zhao et al. (2025), where η is the KL-regularization parameter. A direct implication of this result is that by choosing η to be sufficiently small, one can always guarantee an arbitrarily small KL-regularized regret. Indeed, a small η indicates that the KL-regularized optimal policy π_{η}^{\star} remains very close to the reference policy π_{ref} , which makes this result appear reasonable. However, when $\pi_{\eta}^{\star} \approx \pi_{\text{ref}}$, the learner gains little to no improvement, which is undesirable since the goal is to discover a strictly better policy than the reference policy. To address this, we also consider the notion of *unregularized regret* (Equation 3), as in standard bandit settings. This regret can be large when the policy remains close to π_{ref} (i.e., for small η) but far from the *unregularized optimal policy* π^{\star} . Minimizing the unregularized regret allows us to directly pursue the unregularized optimal policy π^{\star} , rather than being limited to the KL-regularized solution π_{η}^{\star} . This naturally raises the hypothesis that η should be chosen carefully to minimize the unregularized regret, which leads to our second question:

By choosing η appropriately, can we achieve sublinear unregularized regret, still without additional exploration techniques?

In this paper, we answer these questions affirmatively. We begin by analyzing the KL-regularized (adversarial) contextual bandit setting and then extend our analysis to RLHF. To consider general algorithms, we assume access to an online regression oracle (Foster & Rakhlin, 2020), while the offline regression oracle is discussed in Appendix F. Our main contributions are summarized as:

- KL-regularized regret. In KL-regularized contextual bandits, we establish a KL-regularized regret bound of $\mathcal{O}(\eta \mathrm{Reg}_{\mathrm{Sq}}(T) + \eta \log(1/\delta))$, where η is the regularization parameter, $\mathrm{Reg}_{\mathrm{Sq}}(T)$ is the online regression oracle bound, and δ is the failure probability (Theorem 1). This result is achieved solely through KL-regularization, without relying on any additional exploration techniques. To our best knowledge, this is the first result to show the provable efficiency of the KL-regularization-only approach. Since $\mathrm{Reg}_{\mathrm{Sq}}(T) = \mathcal{O}(\log T)$ can be attained by suitable regression oracles for a wide range of reward functions—including linear, generalized linear, and bounded eluder-dimension function classes—we achieve logarithmic KL-regularized regret.
- Unregularized regret. By setting $\eta = \Theta\left(\sqrt{DT/(\mathrm{Reg_{Sq}}(T) + \log \delta^{-1})}\right)$, we obtain an unregularized regret of $\mathcal{O}(\sqrt{DT(\mathrm{Reg_{Sq}}(T) + \log \delta^{-1})})$, where $D = \frac{1}{T}\sum_{t=1}^{T}\mathrm{KL}\left(\pi^{\star}(\cdot\|x_t)\|\pi_{\mathrm{ref}}(\cdot\|x_t)\right)$ (Theorem 1). To the best of our knowledge, this is the first unregularized regret bound for KL-regularized contextual bandits attained solely through KL-regularization-induced exploration.
- $\tilde{\mathcal{O}}(\sqrt{dT\log N})$ regret in linear contextual bandits. By setting π_{ref} to the uniform random policy and choosing $\eta = \Theta(\sqrt{T\log N/(d\log T)})$, we obtain an (unregularized) regret bound of $\tilde{\mathcal{O}}(\sqrt{dT\log N})$ for linear contextual bandits (Theorem 2), where d is the feature dimension, N is the number of arms, and $\tilde{\mathcal{O}}(\cdot)$ hides constant and logarithmic factors (excluding those in N). To the best of our knowledge, this is the first $\tilde{\mathcal{O}}(\sqrt{dT\log N})$ -type regret achieved without using on supLin-type algorithms (Auer, 2002; Chu et al., 2011; Li et al., 2019), which are known to be impractical. Hence, this is the first practical algorithm to achieve minimax optimal regret for finite-armed linear contextual bandits.
- Extension to RLHF. We further establish similar regret bounds in the RLHF setting, with only an additional factor due to the non-linearity of the Bradley–Terry model (Theorems 3 and E.1).

2 RELATED WORKS

Online RLHF. Early works in online RLHF trace back to the dueling bandits literature (Yue et al., 2012; Zoghi et al., 2015; Saha & Gopalan, 2018; Bengs et al., 2021) and were later extended to the reinforcement learning setting (Xu et al., 2020; Novoseller et al., 2020; Chen et al., 2022; Saha et al., 2023; Zhan et al., 2023b; Wu & Sun, 2023). More recently, Xiong et al. (2023); Ye et al. (2024) introduced provably efficient algorithms under the KL-regularized objective using UCB-style exploration. These were further refined by methods that employ optimistically biased optimization targets (Xie et al., 2024; Liu et al., 2024; Cen et al., 2024). The most closely related works are Zhao et al. (2024; 2025), which also study the KL-regularized objective and establish $\mathcal{O}(\eta \log T)$ KL-regularized regret (or $\mathcal{O}(\eta/\epsilon)$ suboptimality gap). However, all of these prior approaches depend on

additional exploration mechanisms. In contrast, our work demonstrates—for the first time, to the best of our knowledge—that KL-regularization alone suffices to achieve sublinear regret in both the regularized and unregularized forms. For additional related work, see Appendix A.

3 PROBLEM SETUP

Notations. Given a set \mathcal{X} , we use $|\mathcal{X}|$ to denote its cardinality. For a positive integer, n, we denote $[n] := \{1, 2, \dots, n\}$. Let N denote the size of the action space. We write $\mathcal{O}(\cdot)$ for asymptotics up to constants and $\tilde{\mathcal{O}}(\cdot)$ when also hiding logarithmic factors (except in N). For a function class \mathcal{F} , we denote by $\mathcal{N}_{\mathcal{F}}(\epsilon)$ its ϵ -covering number.

3.1 KL-REGULARIZED CONTEXTUAL BANDITS

In the KL-regularized contextual bandits, at each round $t \in [T]$, the learner observes a context $x_t \in \mathcal{X}$ (which may be provided *adversarially*) and then selects an action $a_t \in \mathcal{A}$, where \mathcal{X} is the context space and \mathcal{A} is the action space. The learner then receives a reward $r_t \in [0, 1]$, given by:

$$r_t = R^{\star}(x_t, a_t) + \epsilon_t,$$

where $R^{\star}(x_t, a_t)$ is the unknown expected reward function, and ϵ_t is independent, zero-mean, and 1-sub-Gaussian. In this paper, we consider a general reward function class $\mathcal{R} \subseteq \{R : \mathcal{X} \times \mathcal{A} \to [0, 1]\}$, which can be a class of parametric functions, nonparametric functions, neural networks, etc. We assume the standard *realizability* (Chu et al., 2011; Agarwal et al., 2012; Foster et al., 2018a; Foster & Rakhlin, 2020; Simchi-Levi & Xu, 2022) throughout the paper, including the RLHF setting.

Assumption 1 (Realizability). The true reward function is contained in \mathcal{R} , i.e., $R^* \in \mathcal{R}$.

KL-Regularized Objective. We consider a *KL-regularized* reward objective, defined for a regularization parameter $\eta > 0$, as:

$$J_t^{\eta}(\pi, R) := \mathbb{E}_{a \sim \pi(\cdot | x_t)} \left[R(x_t, a) \right] - \frac{1}{\eta} \operatorname{KL} \left(\pi(\cdot | x_t) \| \pi_{\operatorname{ref}}(\cdot | x_t) \right), \quad \forall t \geqslant 1, \tag{1}$$

where π_{ref} is the reference policy known to the learner. When π_{ref} is uniform, Equation 1 reduces to the entropy-regularized objective that encourages diverse actions and enhances robustness (Williams, 1992; Levine & Koltun, 2013; Levine et al., 2016; Haarnoja et al., 2018), which is also closely-related to the generative flow networks (GFlowNets) (Bengio et al., 2021; 2023; Tiapkin et al., 2024). When π_{ref} is instead chosen as a base model, KL regularization has been widely adopted for RL fine-tuning of large language models (Ouyang et al., 2022; Rafailov et al., 2023). It has also been studied in online learning (Cai et al., 2020; He et al., 2022) and convex optimization (Neu et al., 2017).

Following prior work (Peters & Schaal, 2007; Rafailov et al., 2023; Zhang, 2023), it is straightforward to show that the optimal solution to the objective in Equation 1 has the following form:

$$\pi_R^{\eta}(a|x) = \frac{1}{Z_R(x)} \pi_{\text{ref}}(a|x) \exp(\eta R(x, a)), \qquad (2)$$

where $Z_R(x) := \mathbb{E}_{a \sim \pi_{\text{ref}}(\cdot|x)} \exp(\eta R(x,a))$ is the normalization constant. A full derivation can be found in Appendix A.1 of Rafailov et al. (2023).

3.2 REINFORCEMENT LEARNING FROM HUMAN FEEDBACK (RLHF)

In the RLHF problem (Ouyang et al., 2022)—more specifically, the contextual *dueling* bandit problem with a KL-regularized objective—the learner at each round $t \in [T]$ observes a context $x_t \in \mathcal{X}$ (possibly provided *adversarially*) and selects two actions $a_t^1, a_t^2 \in \mathcal{A}$, where \mathcal{X} is the context space and \mathcal{A} the action space. The learner then receives relative preference feedback between the two actions, rather than a scalar reward. In this paper, we consider the Bradley-Terry Model (Bradley & Terry, 1952), where the probability of a^1 is preferred over a^2 (denoted by $a^1 > a^2$) is given by

$$\mathbb{P}(a^1 > a^2 | x, a^1, a^2) = \sigma(R^*(x, a^1) - R^*(x, a^2)),$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function, and $R^\star: \mathcal{X} \times \mathcal{A} \to [0,1]$ the *unknown true* reward function. We denote $\mathcal{R} \subseteq \{R: \mathcal{X} \times \mathcal{A} \to [0,1]\}$ as the class of reward functions. To capture the non-linearity of the sigmoid function, we define $\kappa := \sup_{R \in \mathcal{R}, x \in \mathcal{X}, a \in \mathcal{A}} 1/\dot{\sigma}(R(x,a))$. As in the bandit setting, we update the policy by optimizing the KL-regularized reward objective (Equation 1).

3.3 KL-REGULARIZED AND UNREGULARIZED REGRET

We study two types of regret to more comprehensively evaluate the performance of our algorithm.

KL-regularized regret. Let $\pi_{\eta}^{\star}(\cdot|x_t) = \operatorname{argmax}_{\pi} J_t^{\eta}(\pi, R^{\star})$ denote the *KL-regularized optimal policy*. Our objective is to minimize the cumulative regret, defined as:

$$\mathbf{Regret}_{\mathrm{KL}}(T,\eta) := \sum_{t=1}^{T} \left(J_t^{\eta}(\pi_{\eta}^{\star},R^{\star}) - J_t^{\eta}(\pi_t,R^{\star}) \right).$$

This KL-regularized regret has been extensively studied in the prior literature (Xiong et al., 2023; Ye et al., 2024; Song et al., 2024; Zhao et al., 2024; 2025).

Unregularized regret. Beyond the KL-regularized regret, we also measure performance relative to the unregularized optimal policy $\pi^{\star}(\cdot|x_t) = \arg\max_{\pi} \mathbb{E}_{a \sim \pi(\cdot|x_t)}[R^{\star}(x_t, a)]$, and define the corresponding regret as follows:

$$\mathbf{Regret}(T) := \sum_{t=1}^{T} \left(\mathbb{E}_{a \sim \pi^{\star}(\cdot|x_t)} [R^{\star}(x_t, a)] - \mathbb{E}_{a \sim \pi_t(\cdot|x_t)} [R^{\star}(x_t, a)] \right). \tag{3}$$

The notion of this regret is standard in conventional bandit problems. This metric enables a more direct evaluation of how closely the learned policies approach the unregularized optimal policy.

4 KL-REGULARIZED CONTEXTUAL BANDITS

In this section, we consider KL-regularized contextual bandit problems. We introduce the notion of an online regression oracle (Subsection 4.1), present our algorithm KL-EXP together with its regret bounds (Subsection 4.2), and provide a proof sketch (Subsection 4.3).

4.1 Squared-loss online regression oracle.

We assume access to a squared-loss online regression oracle (Foster & Rakhlin, 2020), denoted by OracleSq. At each round t, OracleSq outputs a reward estimator

$$\widehat{R}_t \leftarrow \text{OracleSq}_t((x_t, a_t); (x_1, a_1, r_1), \dots, (x_{t-1}, a_{t-1}, r_{t-1})). \tag{4}$$

Its prediction error is assumed to be bounded with respect to the true reward R^* .

Assumption 2 (Guarantee of OracleSq). We assume that, for every sequence $x_{1:T}, a_{1:T}, r_{1:T}$, there exists regret bound $Reg_{Sq}(T)$ such that the regression oracle OracleSq satisfies

$$\sum_{t=1}^T \left(\widehat{R}_t(x_t, a_t) - r_t \right)^2 - \sum_{t=1}^T \left(R^{\star}(x_t, a_t) - r_t \right)^2 \leqslant \operatorname{Reg}_{\operatorname{Sq}}(T).$$

An important advantage of Assumption 2 is that it places no restriction on how the estimator \hat{R}_t is obtained; in particular, it does not require solving ERM exactly. Instead, \hat{R}_t can be computed via iterative methods such as (stochastic) gradient descent and implemented in an online or streaming manner, which is crucial for large-scale modern machine learning. Under realizability, Assumption 2 is weaker than Assumption 2a in Foster & Rakhlin (2020), since we compete only against the fixed R^\star , whereas they compete against the best predictor over the sequence.

The online squared-loss regression problem is well studied, with efficient algorithms and regret guarantees for many function classes.

Example 1 (Linear classes). When $R^* \in \mathcal{R}$ and the reward function class \mathcal{R} is linear, i.e., $\mathcal{R} = \{R : R = \phi(x, a)^\top \theta, \theta \in \mathbb{R}^d, \|\theta\|_2 \leq 1\}$, where $\phi(x, a) \in \mathbb{R}^d$ is a known feature map satisfying $\|\phi(x, a)\|_2 \leq 1$, choosing OracleSq as the Vovk–Azoury–Warmuth forecaster (Vovk, 1997; Azoury & Warmuth, 2001) yields $\operatorname{Reg}_{\operatorname{Sq}}(T) = \mathcal{O}(d \log(T/d))$.

Example 2 (Generalized linear models (GLMs)). For a fixed non-decreasing 1-Lipschitz link function $\mu: \mathbb{R} \to [0,1]$, define the reward function class $\mathcal{R} = \{R: R = \mu(\phi(x,a)^\top \theta), \theta \in \mathbb{R}^d, \|\theta\|_2 \leq 1\}$, where $\phi(x,a) \in \mathbb{R}^d$ is a known feature map with $\|\phi(x,a)\|_2 \leq 1$. If $R^\star \in \mathcal{R}$, then the GLMtron algorithm (Kakade et al., 2011) guarantees $\mathrm{Reg}_{\mathrm{Sq}}(T) = \mathcal{O}(\kappa_\mu^2 d \log(T/d))$, where $1/\dot{\mu} \leq \kappa_\mu$.

Algorithm 1 KL-EXP (**KL**-regularized **EXP**onential-weights algorithm)

 \triangleright 0EP0 if RLHF

- 1: **Inputs:** regularization parameter η , reference policy π_{ref} , online regression oracle OracleSq.
- 2: **Initialize:** choose any $\hat{R}_1 \in \mathcal{R}$.
- 3: **for** round t = 1 to T **do**
- 4: Observe context $x_t \in \mathcal{X}$.
- 5: Compute policy $\pi_t(\cdot|x_t) \propto \pi_{\text{ref}}(\cdot|x_t) \exp(\eta \hat{R}_t(x_t,\cdot))$ via Equation 2.
- 6: Sample action $a_t \sim \pi_t(\cdot|x_t)$ and receive reward r_t . \triangleright or $a_t^1, a_t^2 \sim \pi_t(\cdot|x_t)$ if <u>RLHF</u>
- 7: Update R_{t+1} using OracleSq via Equation 4. \triangleright or update using OracleLog if <u>RLHF</u>
- 8: end for

Example 3 (Bounded eluder dimension, Russo & Van Roy, 2013). When $R^* \in \mathcal{R}$ and the reward function class \mathcal{R} has bounded eluder dimension, the empirical risk minimization (ERM) algorithm achieves, with probability at least $1 - \delta$, $\text{Reg}_{S_0}(T) = \mathcal{O}(d_E \log(\mathcal{N}_{\mathcal{R}}(\epsilon)T))$ (Lemma C.2).

For additional examples, the reader is referred to Foster & Rakhlin (2020) for high-dimensional linear models, Banach spaces, and RKHS, and to Deb et al. (2024) for neural networks.

4.2 ALGORITHM AND MAIN RESULTS

We present our KL-regularized EXPonential-weights algorithm, KL-EXP, in Algorithm 1. At each round $t \in [T]$, the algorithm observes the context $x_t \in \mathcal{X}$ and computes the policy π_t by solving the KL-regularized objective in Equation 1, which admits the closed-form solution given in Equation 2. The algorithm then samples an action $a_t \sim \pi_t(\cdot|x_t)$ and receives a reward r_t . Finally, it updates the reward estimator \hat{R}_{t+1} for the next round using the squared-loss online regression oracle (Equation 4). The main guarantees for the algorithm are stated below, with the proof deferred to Appendix B.

Theorem 1 (Regret of KL-EXP). Let $D := \frac{1}{T} \sum_{t=1}^{T} \mathrm{KL} \big(\pi^{\star}(\cdot \| x_t) \| \pi_{\mathrm{ref}}(\cdot \| x_t) \big)$. Under Assumption 2, for any $\delta > 0$, with probability at least $1 - \delta$, KL-EXP (Algorithm 1) guarantees

$$\begin{split} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) &= \mathcal{O}\Big(\eta \mathrm{Reg}_{\mathrm{Sq}}(T) + \eta \log(1/\delta)\Big) \quad \textit{and} \\ \mathbf{Regret}(T) &= \mathcal{O}\bigg(\eta \mathrm{Reg}_{\mathrm{Sq}}(T) + \eta \log(1/\delta) + \frac{DT}{\eta}\bigg) \,. \end{split}$$

Result 1: Logarithmic KL-regularized regret. Theorem 1 shows that the KL-regularized regret of KL-EXP scales with $\operatorname{Reg_{Sq}}(T)$, resulting in logarithmic regret in T across a broad range of function classes. For example, when $\delta = \Theta(T^{-1})$, we obtain $\mathcal{O}(\eta d \log T)$ for linear classes (Example 1), $\mathcal{O}(\eta \kappa_{\mu}^2 d \log T)$ for generalized linear models (Example 2), and $\mathcal{O}(\eta d_{\rm E} \log(\mathcal{N}_{\mathcal{R}}(\epsilon)T))$ for function classes with bounded eluder dimension (Russo & Van Roy, 2013) (Example 3). Hence, Theorem 1 shows that logarithmic KL-regularized regret in T can be achieved without the *auxiliary exploration* methods (e.g., UCB-based strategies). In contrast, prior works such as Xiong et al. (2023; 2024); Xie et al. (2024) obtained $\mathcal{O}(\sqrt{T})$ KL-regularized regret (or $\mathcal{O}(1/\epsilon^2)$ sample complexity), and more recently, Zhao et al. (2024; 2025) established $\mathcal{O}(\eta \log T)$ KL-regularized regret (or $\mathcal{O}(\eta/\epsilon)$ sample complexity), all of which depend on the additional exploration strategies. To the best of our knowledge, this is the first result that achieves logarithmic KL-regularized regret without any additional exploration, highlighting the key insight that the KL-regularized objective alone provides sufficient exploration in contextual dueling bandits and RLHF.

Remark 1 (Comparison with Zhao et al. (2025)). For classes with bounded eluder dimension, we recover the regret bound of Zhao et al. (2025), $\mathcal{O}(\eta d_E \log(\mathcal{N}_{\mathcal{R}}(\epsilon)T))$. Unlike Zhao et al. (2025), however, our algorithm does not require prior knowledge of the eluder dimension (Russo & Van Roy, 2013), which is typically unknown in practice. The full proof is provided in Appendix C.

Result 2: Unregularized regret and its tightness. With the choice of the regularization parameter $\eta = \Theta(\sqrt{DT/(\text{Reg}_{Sq}(T) + \log \delta^{-1})})$, we obtain $\text{Regret}(T) = \mathcal{O}(\sqrt{DT(\text{Reg}_{Sq}(T) + \log \delta^{-1})})$. The result provides an interesting insight: with *appropriately chosen* η , it is possible to achieve a \sqrt{T} -type regret bound even in conventional (unregularized) contextual bandit problems. To the best

of our knowledge, this is the first unregularized regret bound in KL-regularized contextual bandits achieved purely via KL-regularization—induced exploration.

Remark 2 (Ease of implementation and computational efficiency). *KL-EXP is simple and practical:* it admits a closed-form solution (Equation 2) and, unlike prior approaches (Russo & Van Roy, 2013; Jiang et al., 2017; Jin et al., 2021; Zhao et al., 2025), avoids the explicit computation of exploration terms (e.g., UCB), which can be intractable for large models such as transformers. It is also computationally efficient—for instance, in linear contextual bandits (ignoring oracle computations), the per-round cost is only $\mathcal{O}(N)$, compared to $\mathcal{O}(d^2N)$ for LinucB and LintS.

To demonstrate the tightness of our bound, we consider the uniform reference policy $\pi_{\text{ref}} = \text{Unif}(\mathcal{A})$, under which $\text{KL}(\pi \| \pi_{\text{ref}}) \leq \log N$ holds for any policy π . In this setting, for linear (adversarial) contextual bandits, we obtain the first $\tilde{\mathcal{O}}(\sqrt{dT \log N})$ -type regret bound, to the best of our knowledge.

Theorem 2 (Unregularized regret under linear classes). We denote $N = |\mathcal{A}|$. Under the setting of Theorem 1 with Assumption 1, if we set $\pi_{\text{ref}} = \text{Unif}(\mathcal{A})$ and $\eta = \Theta(\sqrt{T \log N/(d \log T)})$, then with probability at least $1 - \frac{1}{T}$, we have $\text{Regret}(T) = \mathcal{O}(\sqrt{dT \log N \log T})$.

Remark 3 (Minimax-optimality under linear classes). The proof of Theorem 2 follows directly from two facts: $\operatorname{Reg}_{\operatorname{Sq}}(T) = \mathcal{O}(d\log(T/d))$ (Example 1) and $\operatorname{KL}(\pi^{\star}\|\pi_{\operatorname{ref}}) \leq \log N$ when $\pi_{\operatorname{ref}} = \operatorname{Unif}(\mathcal{A})$. We highlight that, in linear contextual bandits, our regret bound $\mathcal{O}(\sqrt{dT\log N\log T})$ is minimax-optimal, matching the order previously attained by supLin-type algorithms (Auer, 2002; Chu et al., 2011; Li et al., 2019). To the best of our knowledge, this is the first $\widetilde{\mathcal{O}}(\sqrt{dT\log N})$ -type regret bound for linear (adversarial) contextual bandits that avoids the impractical "layered data partitioning" technique and explicit UCB computations. Moreover, it matches the lower bound $\Omega(\sqrt{dT\log N\log(T/d)})$ (Li et al., 2019) up to logarithmic d factors, underscoring both the statistical and computational efficiency of our approach.

Further examples for specific function classes are provided in Appendix B.4.

4.3 Proof sketch of Theorem 1

1) Second-order regret decomposition. The regret decomposition is similar to the recent work of Zhao et al. (2025), which establishes logarithmic KL-regularized regret. Define the function $f(x,R) := -\frac{1}{\eta} \log Z_R(x) + \mathbb{E}_{\pi_R^{\eta}} \left[R(x,a) - R^{\star}(x,a) \right]$. Since $R^{\star}(x,a) = \frac{1}{\eta} \log \exp \left(\eta R^{\star}(x,a) \right)$, we obtain

$$J_t^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_t^{\eta}(\pi_t, R^{\star}) = \frac{1}{\eta} \log Z_{R^{\star}}(x) - \frac{1}{\eta} \log Z_{\hat{R}_t}(x) + \mathbb{E}_{a \sim \pi_t(\cdot | x_t)} \left[\hat{R}_t(x, a) - R^{\star}(x, a) \right]$$
$$= f(x, \hat{R}_t) - f(x, R^{\star}).$$

In Zhao et al. (2025), the decomposition takes the alternative form $J_t^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_t^{\eta}(\pi_t, R^{\star}) = f(x, \tilde{R}_t) - f(x, R^{\star})$, where $\tilde{R}_t(x, a) := \hat{R}_t(x, a) + b_t(x, a)$ is the UCB. They then apply the mean value theorem to this expression and leverage optimism to bound $f(x, \tilde{R}_t) - f(x, R^{\star})$.

In contrast, our analysis shows that it suffices to work directly with the oracle estimator \hat{R}_t . Instead of invoking the mean value theorem, we use the exact second-order Taylor expansion of f.

$$f(x, \hat{R}_{t}) - f(x, R^{\star}) = \sum_{a \in \mathcal{A}} \underbrace{\frac{\partial f(x, R^{\star})}{\partial R(x, a)}}_{=0} \Delta R_{t}(x, a)$$

$$+ \int_{0}^{1} (1 - \alpha) \left[\sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} \Delta R_{t}(x, a) \frac{\partial^{2} f(x, R^{\star} + \alpha \Delta R_{t})}{\partial R(x, a') \partial R(x, a)} \Delta R_{t}(x, a') \right] d\alpha$$

$$\leq \eta \mathbb{E}_{a \sim \pi_{t}(\cdot | x_{t})} \left[\left(\hat{R}_{t}(x_{t}, a) - R^{\star}(x_{t}, a) \right)^{2} \right], \tag{5}$$

where $\Delta R_t = \hat{R}_t - R^*$. Note that in the equation, $\frac{\partial f(x,R^*)}{\partial R(x,a)} = 0$, which is one of our key theoretical findings. This result shows that it is unnecessary to rely on optimistic estimators such as UCB. The remaining steps then follow directly from straightforward calculus (see Lemma B.2 for details).

2) Conversion to regression oracle bound. By summing over $t \in [T]$ in Equation 5 and applying Freedman's inequality together with Lemma 4 of Foster & Rakhlin 2020, we obtain

$$\mathbf{Regret}_{\mathrm{KL}}(T,\eta) \leqslant \eta \sum_{t=1}^{T} \mathbb{E}_{a_t \sim \pi_t(\cdot \mid x_t)} \left[\left(\widehat{R}_t(x_t, a_t) - R^{\star}(x_t, a_t) \right)^2 \right] \leqslant 2\eta \mathrm{Reg}_{\mathrm{Sq}}(T) + 16\eta \log \frac{1}{\delta}.$$

This completes the proof of the KL-regularized regret bound.

3) Unregularized regret bound. From the definitions of J_t^{η} and π_{η}^{\star} , together with the non-negativity of the KL divergence, we can bound the unregularized regret as follows:

$$\begin{aligned} \mathbf{Regret}(T) &= \mathbb{E}_{a \sim \pi^{\star}(\cdot \mid x_{t})}[R^{\star}(x_{t}, a)] - \mathbb{E}_{a \sim \pi_{t}(\cdot \mid x_{t})}[R^{\star}(x_{t}, a)] \\ &\leq J_{t}^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_{t}^{\eta}(\pi_{t}, R^{\star}) + \frac{1}{\eta} \operatorname{KL} \left(\pi^{\star}(\cdot \parallel x_{t}) \parallel \pi_{\operatorname{ref}}(\cdot \parallel x_{t})\right). \end{aligned}$$

Summing over $t \in [T]$ and applying the KL-regularized regret bound established above, we complete the proof of Theorem 1.

5 REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

5.1 Log-loss online regression oracle.

Similar to the KL-regularized contextual bandit setting, we assume access to a log-loss online regression oracle (Foster & Krishnamurthy, 2021), denoted by OracleLog. First, we define the binary logarithmic/cross-entropy loss function ("log-loss") at round t as

$$\ell_t(R) := - \left[y_t \log \sigma \left(R(x_t, a_t^1) - R(x_t, a_t^2) \right) + (1 - y_t) \log \sigma \left(R(x_t, a_t^2) - R(x_t, a_t^1) \right) \right], \quad (6)$$

where y_t denote the binary preference label, where $y_t=1$ if a_t^1 is preferred over a_t^2 (i.e., $a_t^1>a_t^2$) and $y_t=0$ otherwise At each round t, OracleLog returns

$$\widehat{R}_t \leftarrow \texttt{OracleLog}_t((x_t, a_t^1, a_t^2); (x_1, a_1^1, a_1^2, y_1), \dots, (x_{t-1}, a_{t-1}^1, a_{t-1}^2, y_{t-1})). \tag{7}$$

Analogous to Assumption 2, we assume that the prediction error of OracleLog is bounded as follows: **Assumption 3** (Guarantee of log-loss regression oracle). We assume that, for every (possibly adaptively chosen) sequence $x_{1:T}$, $a_{1:T}^1$, $a_{1:T}^2$, $y_{1:T}$, there exists regret bound $\operatorname{Reg}_{\operatorname{Log}}(T)$ such that the regression oracle OracleLog satisfies

$$\sum_{t=1}^{T} \ell_t(\hat{R}_t) - \sum_{t=1}^{T} \ell_t(R^*) \leqslant \operatorname{Reg}_{\operatorname{Log}}(T).$$

Example 4 (Linear classes under log-loss). When $R^* \in \mathcal{R}$ and the reward function class \mathcal{R} is linear, we can use the algorithm from Foster et al. (2018b) to obtain $\operatorname{Reg}_{\operatorname{Log}}(T) = \mathcal{O}(d \log(T/d))$.

Similar guarantees are available for kernels, generalized linear models, and many other nonparametric classes, as in the case of the squared-loss online regression oracle (Foster & Krishnamurthy, 2021).

5.2 ALGORITHM AND MAIN RESULTS

We now introduce an algorithm for RLHF problems, OEPO, outlined in the comments of Algorithm 1 and fully described in Algorithm D.1. The overall flow is similar to KL-EXP; however, at each round $t \in [T]$, the current policy samples two actions, $a_t^1, a_t^2 \sim \pi_t(\cdot|x_t)$, and receives preference feedback between them. Another key difference is that the reward estimator \hat{R}_{t+1} is updated using the log-loss online regression oracle OracleLog (Equation 7). When OracleLog is implemented with a gradient-based method (e.g., SGD or Adam), OEPO recovers the practical online RLHF algorithm.

The regret guarantees for OEPO are presented below, with the proofs deferred to Appendix D.

Theorem 3 (Regret of 0EP0). Let $D := \frac{1}{T} \sum_{t=1}^{T} \mathrm{KL} \left(\pi^{\star}(\cdot \| x_t) \| \pi_{\mathrm{ref}}(\cdot \| x_t) \right)$ and $\kappa := \sup_{R,x,a} \frac{1}{\dot{\sigma}(R(x,a))}$. Under Assumption 2, for $\delta > 0$, with probability at least $1 - \delta$, 0EP0 ensures

$$\begin{split} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) &= \mathcal{O}\Big(\eta\kappa^2\mathrm{Reg}_{\mathrm{Log}}(T) + \eta\kappa^2\log(1/\delta)\Big) \quad \textit{and} \\ \mathbf{Regret}(T) &= \mathcal{O}\bigg(\eta\kappa^2\mathrm{Reg}_{\mathrm{Log}}(T) + \eta\kappa^2\log(1/\delta) + \frac{DT}{\eta}\bigg) \,. \end{split}$$

Discussion of Theorem 3. We obtain regret bounds comparable to Theorem 1, up to a κ factor (and differences in oracle prediction error). Such κ -dependence is standard and largely unavoidable in RLHF and dueling bandits (Saha, 2021; Saha et al., 2023; Zhu et al., 2023; Xiong et al., 2023; Zhan et al., 2023b; Das et al., 2024; Xie et al., 2024; Zhao et al., 2024). With the choices $\eta = \Theta(\sqrt{DT/(\kappa^2 \text{Reg}_{\text{Log}}(T))})$ and $\pi_{\text{ref}} = \text{Unif}(\mathcal{A})$, 0EPO achieves unregularized regret $\text{Regret}(T) = \mathcal{O}(\kappa\sqrt{DT}\text{Reg}_{\text{Log}}(T))$. As in Theorem 1, this yields $\tilde{\mathcal{O}}(\sqrt{T})$ regret guarantees for a broad range of function classes (see Foster & Krishnamurthy (2021) for bound on $\text{Reg}_{\text{Log}}(T)$).

Remark 4 (Extension to DPO, Rafailov et al., 2023). *The DPO-variant algorithm (Algorithm E.1)* achieves the same-order regrets, up to differences in the oracle's prediction error (see Appendix E).

6 EXPERIMENTS

6.1 LINEAR CONTEXTUAL BANDITS

In the linear bandit experiments, we consider linear reward function class, i.e., $\mathcal{R} = \{R : R = \phi(x,a)^{\top}\theta, \theta \in \mathbb{R}^d, \|\theta\|_2 \leq 1\}$. For each instance we sample the true parameter $\theta^{\star} \sim \mathcal{N}(0,I_d)$ and normalize it so that $\|\theta^{\star}\|_2 \leq 1$. At each round t, a context $x_t \in \mathcal{X}$ is drawn uniformly at random, with feature vector $\phi(x_t,a) \in \mathbb{R}^d$ lying in the unit ball. We set $d \in \{10,20\}$ and $N = |\mathcal{A}| \in \{50,100\}$. We report cumulative regret averaged over 20 runs, with standard errors.

We compare the performance of our algorithm KL-EXP against four baselines: (i) LinUCB (Li et al., 2010), (ii) LinTS (Agrawal & Goyal, 2013), (iii) LinPHE (Kveton et al., 2020), and (iv) SupLinUCB (Chu et al., 2011). For all baselines, we set the confidence bound to $\sqrt{d \log T}$. For KL-EXP, we use $\pi_{\rm ref} = {\rm Unif}(\mathcal{A})$ and choose $\eta = \sqrt{T}$, which yields the $\tilde{\mathcal{O}}(\sqrt{T})$ unregularized regret. Figure 1 shows that our algorithm consistently and significantly outperforms the baselines across varying d and N, while also achieving faster per-round computation than the others (see Table H.1).

6.2 NEURAL CONTEXTUAL BANDITS

In the neural bandit experiments, we use the neural network reward class \mathcal{R} , instantiated as a two-layer network with input dimension 80 and hidden width 100, equipped with ReLU activations. We evaluate four types of true reward functions: (i) linear: $R^*(x,a) = \phi(x,a)^\top \theta^*$, (ii) quadratic: $R^*(x,a) = (\phi(x,a)^\top \theta^*)^2$, (iii) cosine: $R^*(x,a) = \cos(\pi \phi(x,a)^\top \theta^*)$, and (iv) neural network: $R^* \in \mathcal{R}$. Training is performed with squared loss via SGD (batch size 100, learning rate 0.005). We set N=20, and report cumulative regret averaged over 10 runs with standard errors.

We compare our algorithm KL-EXP against two baselines: (i) NeuralUCB (Zhou et al., 2020) and (ii) NeuralTS (Zhang et al., 2020). For the baselines, we tune the confidence bounds via grid search over $\{1.0, 5.0, 10.0\}$. For KL-EXP, we tune η using grid search over $\{50, 100, 500\}$, and adopt the uniform random reference policy. Figure 2 shows that our algorithm outperforms the baselines across diverse reward structures while running about $10\times$ faster (see Table H.2).

6.3 LLM FINE-TUNING WITH RLHF

In this subsection, we validate our key theoretical insight in the LLM fine-tuning task: properly tuning the regularization parameter η alone is sufficient to induce exploration. Our DPO-variant algorithm, ODPO, coincides with OnlineDPO (Guo et al., 2024) when the regression oracle OracleDPO (defined in Equation E.2) is instantiated using the original DPO optimizer settings (optimizer, batch size, learning rate, and training steps). Since we adopt these original settings, we report the algorithm as OnlineDPO (in Figure 3) rather than ODPO, to avoid confusion.

For experimental details, we follow the iterative DPO pipeline (Xu et al., 2023; Tran et al., 2023; Dong et al., 2024; Xie et al., 2024) from Dong et al. (2024), running T=3 total iterations with large batches of pairs sampled from π_t . We use the same base model (Llama-3-8B-Flow-SFT¹), prompt sets for each iteration², and true preference model for generating feedback³ as in Dong et al. (2024);

¹https://huggingface.co/RLHFlow/LLaMA3-SFT

²https://huggingface.co/datasets/RLHFlow/iterative-prompt-v1-iter2-20K

³https://huggingface.co/RLHFlow/pair-preference-model-LLaMA3-8B

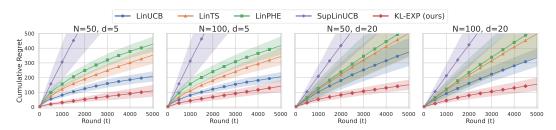


Figure 1: Cumulative regret in linear bandits with $d \in \{10, 20\}$ and $N = |\mathcal{A}| \in \{50, 100\}$.

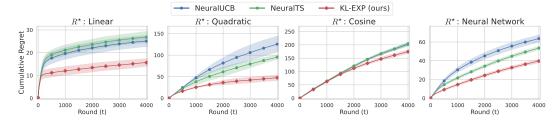


Figure 2: Cumulative regret in neural bandits under different true reward functions.

	Llama-3-8B-Flow	Llama-3-8B-Flow	XPO			OnlineDPO $(\eta=12.5)$		
	-SFT	-Final	iter1	iter2	iter3	iter1	iter2	iter3
Accuracy (%)	59.11	60.47	61.03	61.23	61.55	61.49	61.83	62.09

Figure 3: Average accuracy (%) over 17 benchmarks: OnlineDPO ($\eta = 12.5$, a special case of ODPO) vs. existing algorithms that use additional exploration techniques.

Xie et al. (2024), ensuring our results are directly comparable to theirs. Across all three iterations, we fix the reference policy π_{ref} to the base model Llama-3-8B-Flow-SFT.

We consider three baselines: (i) Llama-3-8B-Flow-SFT, the reference model; (ii) Llama-3-8B-Flow-Final, the final model from Dong et al. (2024), released on Hugging Face⁴; and (iii) XPO (Xie et al., 2024). To induce exploration, Llama-3-8B-Flow-Final constructs preference pairs by maximizing heuristic uncertainty, while XPO augments the DPO objective with an additional exploration term that encourages the policy to behave optimistically. We evaluate all algorithms on 17 academic and chat benchmarks (Zhong et al., 2023; Nie et al., 2019; Hendrycks et al., 2020; Cobbe et al., 2021; Rein et al., 2024; Chen et al., 2021; Zellers et al., 2019; Sakaguchi et al., 2021; Clark et al., 2018; Lin et al., 2021; Mihaylov et al., 2018; Zellers et al., 2018; Sap et al., 2019; Pilehvar & Camacho-Collados, 2018; Levesque et al., 2012; Socher et al., 2013) and report their average accuracies. Table 3 shows that with a properly chosen $\eta=12.5$, OnlineDPO (or ODPO) outperforms other baseline algorithms that rely on auxiliary exploration methods. This supports our main theoretical finding that additional exploration techniques are not essential in online RLHF. See Appendix H.3 for additional experimental details, per-benchmark results, and results with other values of η .

7 Conclusion

We show, for the first time to our knowledge, that KL-regularization alone is sufficient for achieving sublinear regrets. In particular, the KL-regularized regret scales with the regression oracle bound, which can be logarithmic in T for many function classes. Moreover, by carefully choosing the regularization parameter η , we achieve $\tilde{\mathcal{O}}(\sqrt{T})$ unregularized regret, demonstrating that the policy can be improved beyond the KL-regularized optimum. This highlights the pivotal role of η in attaining sublinear unregularized regret. We leave further refinements of η , such as time-varying schedules, as an important direction for future work.

⁴https://huggingface.co/RLHFlow/LLaMA3-iterative-DPO-final

REFERENCES

- Alekh Agarwal, Miroslav Dudík, Satyen Kale, John Langford, and Robert Schapire. Contextual bandit learning with predictable rewards. In *Artificial Intelligence and Statistics*, pp. 19–26. PMLR, 2012.
- Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pp. 127–135. PMLR, 2013.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Katy S Azoury and Manfred K Warmuth. Relative loss bounds for on-line density estimation with the exponential family of distributions. *Machine learning*, 43(3):211–246, 2001.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Akshay Balsubramani, Zohar Karnin, Robert E Schapire, and Masrour Zoghi. Instance-dependent regret bounds for dueling bandits. In *Conference on Learning Theory*, pp. 336–360. PMLR, 2016.
- Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow network based generative models for non-iterative diverse candidate generation. *Advances in neural information processing systems*, 34:27381–27394, 2021.
- Yoshua Bengio, Salem Lahlou, Tristan Deleu, Edward J Hu, Mo Tiwari, and Emmanuel Bengio. Gflownet foundations. *Journal of Machine Learning Research*, 24(210):1–55, 2023.
- Viktor Bengs, Róbert Busa-Fekete, Adil El Mesaoudi-Paul, and Eyke Hüllermeier. Preference-based online learning with dueling bandits: A survey. *Journal of Machine Learning Research*, 22(7): 1–108, 2021.
- Viktor Bengs, Aadirupa Saha, and Eyke Hüllermeier. Stochastic contextual dueling bandits under linear stochastic transitivity models. In *International Conference on Machine Learning*, pp. 1764–1786. PMLR, 2022.
- Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. In *International Conference on Machine Learning*, pp. 1283–1294. PMLR, 2020.
- Shicong Cen, Jincheng Mei, Katayoon Goshvadi, Hanjun Dai, Tong Yang, Sherry Yang, Dale Schuurmans, Yuejie Chi, and Bo Dai. Value-incentivized preference optimization: A unified approach to online and offline rlhf. *arXiv preprint arXiv:2405.19320*, 2024.
- Jonathan D Chang, Wenhao Zhan, Owen Oertell, Kianté Brantley, Dipendra Misra, Jason D Lee, and Wen Sun. Dataset reset policy optimization for rlhf. *arXiv preprint arXiv:2404.08495*, 2024.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Xiaoyu Chen, Han Zhong, Zhuoran Yang, Zhaoran Wang, and Liwei Wang. Human-in-the-loop: Provably efficient preference-based reinforcement learning with general function approximation. In *International Conference on Machine Learning*, pp. 3773–3793. PMLR, 2022.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 208–214. JMLR Workshop and Conference Proceedings, 2011.

- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv preprint arXiv:1803.05457, 2018.
 - Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
 - Nirjhar Das, Souradip Chakraborty, Aldo Pacchiano, and Sayak Ray Chowdhury. Active preference optimization for sample efficient rlhf. *arXiv preprint arXiv:2402.10500*, 2024.
 - Rohan Deb, Yikun Ban, Shiliang Zuo, Jingrui He, and Arindam Banerjee. Contextual bandits with online neural regression. In *The Twelfth International Conference on Learning Representations*, 2024
 - Qiwei Di, Tao Jin, Yue Wu, Heyang Zhao, Farzad Farnoud, and Quanquan Gu. Variance-aware regret bounds for stochastic contextual dueling bandits. *arXiv preprint arXiv:2310.00968*, 2023.
 - Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf. arXiv preprint arXiv:2405.07863, 2024.
 - Miroslav Dudík, Katja Hofmann, Robert E Schapire, Aleksandrs Slivkins, and Masrour Zoghi. Contextual dueling bandits. In *Conference on Learning Theory*, pp. 563–587. PMLR, 2015.
 - Moein Falahatgar, Yi Hao, Alon Orlitsky, Venkatadheeraj Pichapati, and Vaishakh Ravindrakumar. Maxing and ranking with few assumptions. *Advances in Neural Information Processing Systems*, 30, 2017.
 - Max H Farrell, Tengyuan Liang, and Sanjog Misra. Deep neural networks for estimation and inference. *Econometrica*, 89(1):181–213, 2021.
 - Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with regression oracles. In *International conference on machine learning*, pp. 3199–3210. PMLR, 2020.
 - Dylan Foster, Alekh Agarwal, Miroslav Dudík, Haipeng Luo, and Robert Schapire. Practical contextual bandits with regression oracles. In *International Conference on Machine Learning*, pp. 1539–1548. PMLR, 2018a.
 - Dylan J Foster and Akshay Krishnamurthy. Efficient first-order contextual bandits: Prediction, allocation, and triangular discrimination. *Advances in Neural Information Processing Systems*, 34: 18907–18919, 2021.
 - Dylan J Foster, Satyen Kale, Haipeng Luo, Mehryar Mohri, and Karthik Sridharan. Logistic regression: The importance of being improper. In *Conference on learning theory*, pp. 167–208. PMLR, 2018b.
 - Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making. *arXiv* preprint arXiv:2112.13487, 2021.
 - David A Freedman. On tail probabilities for martingales. *the Annals of Probability*, pp. 100–118, 1975.
 - Zhaolin Gao, Jonathan Chang, Wenhao Zhan, Owen Oertell, Gokul Swamy, Kianté Brantley, Thorsten Joachims, Drew Bagnell, Jason D Lee, and Wen Sun. Rebel: Reinforcement learning via regressing relative rewards. *Advances in Neural Information Processing Systems*, 37:52354–52400, 2024.
 - Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.
 - Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. Pmlr, 2018.

- Jiafan He, Dongruo Zhou, and Quanquan Gu. Near-optimal policy optimization algorithms for learning adversarial linear mixture mdps. In *International Conference on Artificial Intelligence and Statistics*, pp. 4259–4280. PMLR, 2022.
- Reinhard Heckel, Max Simchowitz, Kannan Ramchandran, and Martin Wainwright. Approximate ranking from pairwise comparisons. In *International Conference on Artificial Intelligence and Statistics*, pp. 1057–1066. PMLR, 2018.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Kevin Jamieson, Sumeet Katariya, Atul Deshpande, and Robert Nowak. Sparse dueling bandits. In *Artificial Intelligence and Statistics*, pp. 416–424. PMLR, 2015.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*, pp. 1704–1713. PMLR, 2017.
- Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. *Advances in neural information processing systems*, 34:13406–13418, 2021.
- Sham M Kakade, Varun Kanade, Ohad Shamir, and Adam Kalai. Efficient learning of generalized linear and single index models with isotonic regression. *Advances in Neural Information Processing Systems*, 24, 2011.
- Junpei Komiyama, Junya Honda, Hisashi Kashima, and Hiroshi Nakagawa. Regret lower bound and optimal algorithm in dueling bandit problem. In *Conference on learning theory*, pp. 1141–1154. PMLR, 2015.
- Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Copeland dueling bandit problem: Regret lower bound, optimal algorithm, and computationally efficient algorithm. In *International Conference on Machine Learning*, pp. 1235–1244. PMLR, 2016.
- Branislav Kveton, Csaba Szepesvári, Mohammad Ghavamzadeh, and Craig Boutilier. Perturbedhistory exploration in stochastic linear bandits. In *Uncertainty in Artificial Intelligence*, pp. 530–540. PMLR, 2020.
- John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. *Advances in neural information processing systems*, 20, 2007.
- Junghyun Lee, Se-Young Yun, and Kwang-Sung Jun. A unified confidence sequence for generalized linear models, with applications to bandits. *arXiv preprint arXiv:2407.13977*, 2024.
- Hector J Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. *KR*, 2012(13th):3, 2012.
- Sergey Levine and Vladlen Koltun. Guided policy search. In *International conference on machine learning*, pp. 1–9. PMLR, 2013.
- Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, 2016.
- Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pp. 661–670, 2010.
- Lihong Li, Yu Lu, and Dengyong Zhou. Provably optimal algorithms for generalized linear contextual bandits. In *International Conference on Machine Learning*, pp. 2071–2080. PMLR, 2017.
- Yingkai Li, Yining Wang, and Yuan Zhou. Nearly minimax-optimal regret for linearly parameterized bandits. In *Conference on Learning Theory*, pp. 2173–2174. PMLR, 2019.

- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
 - Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. Provably mitigating overoptimization in RLHF: Your SFT loss is implicitly an adversarial regularizer. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=2cQ31Phke0.
 - Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
 - Gergely Neu, Anders Jonsson, and Vicenç Gómez. A unified view of entropy-regularized markov decision processes. *arXiv preprint arXiv:1705.07798*, 2017.
 - Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial nli: A new benchmark for natural language understanding. *arXiv preprint arXiv:1910.14599*, 2019.
 - Ellen Novoseller, Yibing Wei, Yanan Sui, Yisong Yue, and Joel Burdick. Dueling posterior sampling for preference-based reinforcement learning. In *Conference on Uncertainty in Artificial Intelligence*, pp. 1029–1038. PMLR, 2020.
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730–27744, 2022.
 - Jan Peters and Stefan Schaal. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, pp. 745–750, 2007.
 - Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. *arXiv preprint arXiv:1808.09121*, 2018.
 - Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
 - Alexander Rakhlin and Karthik Sridharan. Online non-parametric regression. In *Conference on Learning Theory*, pp. 1232–1264. PMLR, 2014.
 - Siddartha Y Ramamohan, Arun Rajkumar, and Shivani Agarwal. Dueling bandits: Beyond condorcet winners to general tournament solutions. Advances in Neural Information Processing Systems, 29, 2016.
 - David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
 - Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, pp. 2256–2264, 2013.
 - Addirupa Saha. Optimal algorithms for stochastic contextual preference bandits. *Advances in Neural Information Processing Systems*, 34:30050–30062, 2021.
 - Aadirupa Saha and Aditya Gopalan. Battle of bandits. In *UAI*, pp. 805–814, 2018.
 - Aadirupa Saha and Akshay Krishnamurthy. Efficient and optimal algorithms for contextual dueling bandits under realizability. In *International Conference on Algorithmic Learning Theory*, pp. 968–994. PMLR, 2022.
 - Aadirupa Saha, Tomer Koren, and Yishay Mansour. Adversarial dueling bandits. In *International Conference on Machine Learning*, pp. 9235–9244. PMLR, 2021.

- Aadirupa Saha, Aldo Pacchiano, and Jonathan Lee. Dueling rl: Reinforcement learning with trajectory preferences. In *International conference on artificial intelligence and statistics*, pp. 6263–6289. PMLR, 2023.
 - Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
 - Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
 - Ayush Sawarni, Nirjhar Das, Siddharth Barman, and Gaurav Sinha. Generalized linear bandits with limited adaptivity. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=FTPDBQuT4G.
 - Ayush Sekhari, Karthik Sridharan, Wen Sun, and Runzhe Wu. Contextual bandits and imitation learning with preference-based active queries. *Advances in Neural Information Processing Systems*, 36:11261–11295, 2023.
 - David Simchi-Levi and Yunzong Xu. Bypassing the monster: A faster and simpler optimal algorithm for contextual bandits under realizability. *Mathematics of Operations Research*, 47(3):1904–1931, 2022.
 - Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1631–1642, 2013.
 - Yuda Song, Gokul Swamy, Aarti Singh, J Bagnell, and Wen Sun. The importance of online data: Understanding preference fine-tuning via coverage. *Advances in Neural Information Processing Systems*, 37:12243–12270, 2024.
 - Daniil Tiapkin, Nikita Morozov, Alexey Naumov, and Dmitry P Vetrov. Generative flow networks as entropy-regularized rl. In *International Conference on Artificial Intelligence and Statistics*, pp. 4213–4221. PMLR, 2024.
 - Hoang Tran, Chris Glaze, and Braden Hancock. Iterative dpo alignment. Technical report, Technical report, Snorkel AI, 2023.
 - Volodya Vovk. Competitive on-line linear regression. *Advances in Neural Information Processing Systems*, 10, 1997.
 - Yuanhao Wang, Qinghua Liu, and Chi Jin. Is rlhf more difficult than standard rl? a theoretical perspective. *Advances in Neural Information Processing Systems*, 36:76006–76032, 2023.
 - Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
 - Huasen Wu and Xin Liu. Double thompson sampling for dueling bandits. *Advances in neural information processing systems*, 29, 2016.
 - Runzhe Wu and Wen Sun. Making rl with preference-based feedback efficient via randomization. *arXiv preprint arXiv:2310.14554*, 2023.
- Yue Wu, Tao Jin, Hao Lou, Farzad Farnoud, and Quanquan Gu. Borda regret minimization for generalized linear dueling bandits. *arXiv preprint arXiv:2303.08816*, 2023.
 - Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q*-approximation for sample-efficient rlhf. *arXiv preprint arXiv:2405.21046*, 2024.
 - Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. *arXiv preprint arXiv:2312.11456*, 2023.

- Wei Xiong, Chengshuai Shi, Jiaming Shen, Aviv Rosenberg, Zhen Qin, Daniele Calandriello, Misha Khalman, Rishabh Joshi, Bilal Piot, Mohammad Saleh, et al. Building math agents with multi-turn iterative preference learning. *arXiv* preprint arXiv:2409.02392, 2024.
 - Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, and Jason Weston. Some things are more cringe than others: Preference optimization with the pairwise cringe loss. *CoRR*, 2023.
 - Yichong Xu, Ruosong Wang, Lin Yang, Aarti Singh, and Artur Dubrawski. Preference-based reinforcement learning with finite-time guarantees. *Advances in Neural Information Processing Systems*, 33:18784–18794, 2020.
 - Chenlu Ye, Wei Xiong, Yuheng Zhang, Nan Jiang, and Tong Zhang. A theoretical analysis of nash learning from human feedback under general kl-regularized preference. *CoRR*, 2024.
 - Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The k-armed dueling bandits problem. *Journal of Computer and System Sciences*, 78(5):1538–1556, 2012.
 - Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. Swag: A large-scale adversarial dataset for grounded commonsense inference. *arXiv preprint arXiv:1808.05326*, 2018.
 - Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
 - Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable offline preference-based reinforcement learning. *arXiv preprint arXiv:2305.14816*, 2023a.
 - Wenhao Zhan, Masatoshi Uehara, Wen Sun, and Jason D Lee. Provable reward-agnostic preference-based reinforcement learning. *arXiv preprint arXiv:2305.18505*, 2023b.
 - Tong Zhang. *Mathematical analysis of machine learning algorithms*. Cambridge University Press, 2023.
 - Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural thompson sampling. *arXiv* preprint arXiv:2010.00827, 2020.
 - Heyang Zhao, Chenlu Ye, Quanquan Gu, and Tong Zhang. Sharp analysis for kl-regularized contextual bandits and rlhf. *arXiv preprint arXiv:2411.04625*, 2024.
 - Heyang Zhao, Chenlu Ye, Wei Xiong, Quanquan Gu, and Tong Zhang. Logarithmic regret for online kl-regularized reinforcement learning. *arXiv preprint arXiv:2502.07460*, 2025.
 - Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv* preprint arXiv:2304.06364, 2023.
 - Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based exploration. In *International Conference on Machine Learning*, pp. 11492–11502. PMLR, 2020.
 - Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In *International Conference on Machine Learning*, pp. 43037–43067. PMLR, 2023.
 - Masrour Zoghi, Shimon Whiteson, Remi Munos, and Maarten Rijke. Relative upper confidence bound for the k-armed dueling bandit problem. In *International conference on machine learning*, pp. 10–18. PMLR, 2014.
 - Masrour Zoghi, Zohar S Karnin, Shimon Whiteson, and Maarten De Rijke. Copeland dueling bandits. *Advances in neural information processing systems*, 28, 2015.

Appendix

Table of Contents

A	Further Related Work							
В	Proof of Theorem 1	17						
	B.1 Main Proof of Theorem 1	17						
	B.2 Proofs of Lemmas for Theorem 1	19						
	B.3 Supporting Results for Lemma B.2	23						
	B.4 Discussion on Specific Function Classes	24						
C	Case: \mathcal{R} with Bounded Eluder Dimension (Remark 1)	25						
D	Proof of Theorem 3	27						
	D.1 Main Proof of Theorem 3	27						
	D.2 Proofs of Lemmas for Theorem 3	28						
E	Extension to Direct Preference Optimization (DPO)	30						
	E.1 Comparison to Lower Bound in Proposition 2.1 of Xie et al. (2024)	31						
F	KL-Regularized Contextual Bandits with Offline Regression Oracle	32						
	F.1 Offline Regression Oracle	32						
	F.2 Algorithm and Results	33						
	F.3 Main Proof of Theorem F.1	33						
G	Technical Lemmas	34						
Н	Additional Experimental Results	35						
	H.1 Computation Cost in Linear Bandits	35						
	H.2 Computation Cost in Neural Bandits	35						
	H.3 RLHF Experiments: Details and Additional Results	35						

A FURTHER RELATED WORK

In this section, we provide additional related work that complements Section 2.

Dueling bandits. The dueling bandit problem, first introduced by Yue et al. (2012), generalizes the classical multi-armed bandit by replacing direct reward observations with pairwise comparisons: in each round t, the learner chooses two arms and only observes which one is preferred. A challenge in this setting is that there may not exist a single arm that dominates all others under arbitrary preference structures. To deal with this, the literature has proposed several notions of "winners," such as the Condorcet winner (Zoghi et al., 2014; Komiyama et al., 2015), Copeland winner (Zoghi et al., 2015; Wu & Liu, 2016; Komiyama et al., 2016), Borda winner (Jamieson et al., 2015; Falahatgar et al., 2017; Heckel et al., 2018; Saha et al., 2021; Wu et al., 2023), and von Neumann winner (Ramamohan et al., 2016; Dudík et al., 2015; Balsubramani et al., 2016), each of which comes with its own performance criterion.

To incorporate contextual information, Saha (2021) introduced the contextual dueling bandit with a Bradley–Terry–Luce (BTL) model (Bradley & Terry, 1952), where pairwise preferences are determined by latent arm rewards. Building on this line, Bengs et al. (2022) analyzed a contextual linear stochastic transitivity model, and Di et al. (2023) proposed a layered algorithm with variance-sensitive regret guarantees.

 Another line of research avoids parametric reward models and instead assumes that preferences are generated by a more general function class. For instance, Saha & Krishnamurthy (2022) developed an algorithm with optimal regret guarantees for K-armed contextual dueling bandits, and Sekhari et al. (2023) further extended the framework with algorithms that provide theoretical guarantees not only on regret but also on query complexity.

However, existing dueling bandit frameworks do not consider the KL-regularized objective, which is the main focus of our work.

RLHF theory. Motivated by the remarkable success of RLHF in fine-tuning LLMs, its theoretical foundations have recently become an active research topic. Much of the existing work focuses on the offline RLHF setting (Zhu et al., 2023; Zhan et al., 2023a), which is complementary to ours. Another line of research studies hybrid RLHF, where offline data are incorporated into an online RL procedure (Xiong et al., 2023; Gao et al., 2024; Chang et al., 2024).

In the context of online RLHF, much of the prior work (Xu et al., 2020; Novoseller et al., 2020; Saha et al., 2023; Xiong et al., 2023; Wu & Sun, 2023) has focused on the special case of tabular MDPs or linear MDPs (or linear reward models when the horizon length is 1), establishing sample complexity or regret bounds in this setting. The exploration bonuses used in these algorithms are specifically designed for linear structures and thus do not extend naturally to the more general function approximation regime we study (e.g., for LLMs).

To go beyond linear models, Chen et al. (2022); Wang et al. (2023); Ye et al. (2024) investigate general function approximation under the assumption of prior knowledge of the eluder dimension (Russo & Van Roy, 2013), which is notoriously difficult to quantify in practice, especially for LLMs. More recently, Zhao et al. (2025) leveraged the properties of KL-regularization to establish the first $\mathcal{O}(\eta \log T)$ KL-regularized regret bound, again assuming prior knowledge of the eluder dimension. These approaches also require solving a complex optimization problem to compute the exploration terms, raising concerns about their practicality for large-scale language models. In parallel, Zhao et al. (2024) achieved a $\mathcal{O}(\eta/\epsilon)$ KL-regularized suboptimality gap by relying on a forced exploration phase, whose length depends on the coverage coefficient—another quantity that is difficult to determine in practice.

To improve practicality under general function approximation, Xie et al. (2024); Liu et al. (2024); Cen et al. (2024) proposed value-incentivized exploration methods that optimize the policy against optimistically biased targets. However, the optimization problems in these approaches do not admit closed-form solutions, and they introduce an additional exploration parameter α that must be tuned, which can make implementation sensitive to hyperparameter choices.

To the best of our knowledge, all existing online RLHF works rely on auxiliary exploration methods beyond KL-regularization. In contrast, our algorithm KL-EXP relies solely on KL-regularization. Moreover, it requires no prior knowledge of any complexity measure, admits a closed-form solution Equation 2, and is thus easy to implement.

B Proof of Theorem 1

In this section, we present the proof of Theorem 1.

B.1 MAIN PROOF OF THEOREM 1

Define $M_t := (\hat{R}_t(x_t, a_t) - r_t)^2 - (R^*(x_t, a_t) - r_t)^2$ and $Z_t := \mathbb{E}[M_t \mid \mathcal{F}_{t-1}] - M_t$, where $\mathcal{F}_{t-1} = \sigma(x_1, a_1, r_1, \dots, x_{t-1}, a_{t-1}, r_{t-1}, x_t)$ is the filtration up to round t-1. The following lemma establishes that these random variables are both bounded and self-bounding.

Lemma B.1 (Lemma 4 of Foster & Rakhlin 2020). Let \mathcal{F}_{t-1} be the filtration up to round t-1, i.e., $\mathcal{F}_{t-1} = \sigma(x_1, a_1, r_1, \dots, x_{t-1}, a_{t-1}, r_{t-1}, x_t)$. Define $M_t := (\hat{R}_t(x_t, a_t) - r_t)^2 - (R^*(x_t, a_t) - r_t)^2$ and $Z_t := \mathbb{E}[M_t \mid \mathcal{F}_{t-1}] - M_t$. Then, the following properties hold:

• $|Z_t| \leq 1$.

•
$$\mathbb{E}[M_t \mid \mathcal{F}_{t-1}] = \mathbb{E}_{a \sim \pi_t(\cdot \mid x_t)} \left[(\widehat{R}_t(x_t, a_t) - R^{\star}(x_t, a_t))^2 \right].$$

•
$$\mathbb{E}[Z_t^2 \mid \mathcal{F}_{t-1}] \leq 4\mathbb{E}[M_t \mid \mathcal{F}_{t-1}].$$

We now present a key lemma that is central to the proof of Theorem 1 and crucial for establishing regret guarantees without any additional exploration.

Lemma B.2 (Second-order regret decomposition). For any $t \in [T]$, we have

$$J_t^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_t^{\eta}(\pi_t, R^{\star}) \leqslant \eta \mathbb{E}_{a \sim \pi_t(\cdot | x_t)} \left[\left(\hat{R}_t(x_t, a) - R^{\star}(x_t, a) \right)^2 \right].$$

The proof is deferred to Appendix B.2.1.

Remark B.1 (Comparison with Zhao et al. (2024)). Unlike Lemma 3.9 of Zhao et al. (2024), which bounds the regret $J_t^{\eta}(\pi_{\eta}^*, R^*) - J_t^{\eta}(\pi_t, R^*)$ in terms of the unknown policy $\pi_{f_{\gamma}}^{\eta}$ (where $f_{\gamma} = \gamma \hat{R}_t + (1 - \gamma)R^*$ for some unknown $\gamma \in (0, 1)$), Lemma B.2 shows that our regret bound depends only on the known current policy π_t . Note that in Zhao et al. (2024), handling the unknown policy $\pi_{f_{\gamma}}^{\eta}$ requires a forced sampling phase, and the minimum number of forced sampling rounds depends on difficult-to-estimate quantities such as the data coverage coefficient (Definition 4.5 therein) and the ϵ -covering number of the reward function class. In contrast, our algorithm does not rely on such quantities.

Remark B.2 (Comparison with Zhao et al. (2025)). Unlike Lemma A.1 of Zhao et al. (2025), Lemma B.2 does not rely on the optimism event. Consequently, our algorithm does not require computing the Upper Confidence Bound (UCB) term, which is generally intractable for general function classes.

Lemma B.3 (Unregularized regret decomposition). For any $t \in [T]$, we have

$$\begin{aligned} \mathbf{Regret}(T) &= \mathbb{E}_{a \sim \pi^{\star}(\cdot \mid x_{t})}[R^{\star}(x_{t}, a)] - \mathbb{E}_{a \sim \pi_{t}(\cdot \mid x_{t})}[R^{\star}(x_{t}, a)] \\ &\leqslant J_{t}^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_{t}^{\eta}(\pi_{t}, R^{\star}) + \frac{1}{\eta} \operatorname{KL} \left(\pi^{\star}(\cdot \parallel x_{t}) \parallel \pi_{\operatorname{ref}}(\cdot \parallel x_{t})\right). \end{aligned}$$

The proof is deferred to Appendix B.2.2.

We are now ready to provide the proof of Theorem 1.

Proof of Theorem 1. By Lemma B.2, we can bound the regret as follows:

$$\mathbf{Regret}_{\mathrm{KL}}(T, \eta) = \sum_{t=1}^{T} J_{t}^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_{t}^{\eta}(\pi_{t}, R^{\star})$$

$$\leq \eta \sum_{t=1}^{T} \mathbb{E}_{a_{t} \sim \pi_{t}(\cdot \mid x_{t})} \left[\left(\widehat{R}_{t}(x_{t}, a_{t}) - R^{\star}(x_{t}, a_{t}) \right)^{2} \right]. \tag{B.1}$$

Let $\mathcal{F}_{t-1} = \sigma(x_1, a_1, r_1, \dots, x_{t-1}, a_{t-1}, r_{t-1}, x_t)$ be the filtration up to round t-1. Define $M_t := (\hat{R}_t(x_t, a_t) - r_t)^2 - (R^*(x_t, a_t) - r_t)^2$ and $Z_t := \mathbb{E}[M_t \mid \mathcal{F}_{t-1}] - M_t$. Then, by applying Freedman's inequality (Lemma G.1) with $\beta = 1/8$, with probability at least $1 - \delta$, we have

$$\sum_{t=1}^{T} \mathbb{E}[M_{t} \mid \mathcal{F}_{t-1}] \leq \sum_{t=1}^{T} M_{t} + \frac{1}{8} \sum_{t=1}^{T} \mathbb{E}[Z_{t}^{2} \mid \mathcal{F}_{t-1}] + 8 \log \frac{1}{\delta}$$

$$= \sum_{t=1}^{T} M_{t} + \frac{1}{2} \sum_{t=1}^{T} \mathbb{E}[M_{t} \mid \mathcal{F}_{t-1}] + 8 \log \frac{1}{\delta} \qquad \text{(Lemma B.1)}$$

$$\leq \text{Reg}_{Sq}(T) + \frac{1}{2} \sum_{t=1}^{T} \mathbb{E}[M_{t} \mid \mathcal{F}_{t-1}] + 8 \log \frac{1}{\delta},$$

where the last inequality holds because

$$\sum_{t=1}^{T} M_t = \sum_{t=1}^{T} (\hat{R}_t(x_t, a_t) - r_t)^2 - \sum_{t=1}^{T} (R^*(x_t, a_t) - r_t)^2 \leqslant \text{Reg}_{Sq}(T).$$
 (Assumption 2)

This directly implies

$$\sum_{t=1}^{T} \mathbb{E}[M_t \mid \mathcal{F}_{t-1}] \le 2 \text{Reg}_{Sq}(T) + 16 \log \frac{1}{\delta}.$$
 (B.2)

Plugging Equation B.2 into Equation B.1, we obtain

$$\begin{aligned} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) &\leqslant \eta \sum_{t=1}^{T} \mathbb{E}_{a_{t} \sim \pi_{t}(\cdot \mid x_{t})} \left[\left(\hat{R}_{t}(x_{t}, a_{t}) - R^{\star}(x_{t}, a_{t}) \right)^{2} \right] \\ &= \eta \sum_{t=1}^{T} \mathbb{E}[M_{t} \mid \mathcal{F}_{t-1}] \\ &\leqslant 2\eta \mathrm{Reg}_{\mathrm{Sq}}(T) + 16\eta \log \frac{1}{\delta}. \end{aligned} \tag{Equation B.2}$$

This concludes the proof of the regret bound for the KL-regularized objective.

We now provide the proof of the unregularized regret bound. By summing over $t \in [T]$ on both sides of the result in Lemma B.3, we directly obtain

$$\begin{split} \mathbf{Regret}(T) &\leqslant \sum_{t=1}^{T} \left(J_{t}^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_{t}^{\eta}(\pi_{t}, R^{\star}) \right) + \frac{1}{\eta} \sum_{t=1}^{T} \mathrm{KL} \big(\pi^{\star}(\cdot \| x_{t}) \| \pi_{\mathrm{ref}}(\cdot \| x_{t}) \big) \\ &= \mathbf{Regret}_{\mathrm{KL}}(T, \eta) + \frac{1}{\eta} \sum_{t=1}^{T} \mathrm{KL} \big(\pi^{\star}(\cdot \| x_{t}) \| \pi_{\mathrm{ref}}(\cdot \| x_{t}) \big) \quad \text{(Definition of } \mathbf{Regret}_{\mathrm{KL}}(T, \eta) \text{)} \\ &= \mathbf{Regret}_{\mathrm{KL}}(T, \eta) + \frac{DT}{\eta} \qquad \qquad (D := \frac{1}{T} \sum_{t=1}^{T} \mathrm{KL} \big(\pi^{\star}(\cdot \| x_{t}) \| \pi_{\mathrm{ref}}(\cdot \| x_{t}) \big) \text{)} \\ &= \mathcal{O} \bigg(\eta \mathrm{Reg}_{\mathrm{Sq}}(T) + \eta \log(1/\delta) + \frac{DT}{\eta} \bigg) \,. \end{split}$$

Hence, the proof of Theorem 1 is complete.

B.2 Proofs of Lemmas for Theorem 1

B.2.1 PROOF OF LEMMA B.2

Proof of Lemma B.2. For simplicity, we use the shorthand $\mathbb{E}_{\pi}[\cdot] = \mathbb{E}_{a \sim \pi(\cdot|x)}[\cdot]$. Noting that $R^{\star}(x, a) = \frac{1}{n} \log \exp(\eta R^{\star}(x, a))$, we have

$$\mathbb{E}_{\pi_{\eta}^{\star}} \left[R^{\star}(x, a) - \frac{1}{\eta} \log \frac{\pi_{\eta}^{\star}(a|x)}{\pi_{\text{ref}}(a|x)} \right] - \mathbb{E}_{\pi_{t}} \left[R^{\star}(x, a) - \frac{1}{\eta} \log \frac{\pi_{t}(a|x)}{\pi_{\text{ref}}(a|x)} \right] \\
= \frac{1}{\eta} \mathbb{E}_{\pi_{\eta}^{\star}} \left[\log \frac{\pi_{\text{ref}}(a|x) \cdot \exp(\eta R^{\star}(x, a))}{\pi_{\eta}^{\star}(a|x)} \right] - \frac{1}{\eta} \mathbb{E}_{\pi_{t}} \left[\log \frac{\pi_{\text{ref}}(a|x) \cdot \exp(\eta R^{\star}(x, a))}{\pi_{t}(a|x)} \right] \\
= \frac{1}{\eta} \mathbb{E}_{\pi_{\eta}^{\star}} \left[\log \frac{\pi_{\text{ref}}(a|x) \cdot \exp(\eta R^{\star}(x, a))}{\pi_{\eta}^{\star}(a|x)} \right] - \frac{1}{\eta} \mathbb{E}_{\pi_{t}} \left[\log \frac{\pi_{\text{ref}}(a|x) \cdot \exp(\eta \hat{R}(x, a))}{\pi_{t}(a|x)} \right] \\
+ \mathbb{E}_{\pi_{t}} \left[\hat{R}_{t}(x, a) - R^{\star}(x, a) \right] \\
= \frac{1}{\eta} \log Z_{R^{\star}}(x) - \frac{1}{\eta} \log Z_{\hat{R}_{t}}(x) + \mathbb{E}_{\pi_{t}} \left[\hat{R}_{t}(x, a) - R^{\star}(x, a) \right], \tag{B.3}$$

where the last equality holds because

$$\frac{\pi_{\text{ref}}(a|x) \cdot \exp\left(\eta R^{\star}(x, a)\right)}{\pi_{n}^{\star}(a|x)} = \frac{\pi_{\text{ref}}(a|x) \cdot \exp\left(\eta R^{\star}(x, a)\right)}{\pi_{\text{ref}}(a|x) \cdot \exp\left(\eta R^{\star}(x, a)\right) / Z_{R^{\star}}(x)} = Z_{R^{\star}}(x),$$

and

$$\frac{\pi_{\text{ref}}(a|x) \cdot \exp\left(\eta R^{\star}(x,a)\right)}{\pi_{t}(a|x)} = \frac{\pi_{\text{ref}}(a|x) \cdot \exp\left(\eta \widehat{R}_{t}(x,a)\right)}{\pi_{\text{ref}}(a|x) \cdot \exp\left(\eta \widehat{R}_{t}(x,a)\right) / Z_{\widehat{R}_{t}}(x)} = Z_{\widehat{R}_{t}}(x),$$

Define the function $f: \mathcal{X} \times \mathcal{R} \to \mathbb{R}$ as follows:

$$f(x,R) := -\frac{1}{\eta} \log Z_R(x) + \sum_{a \in \mathcal{A}} \underbrace{\frac{\pi_{\text{ref}}(a|x) \cdot \exp(\eta R(x,a))}{Z_R(x)}}_{=\pi_R^{\eta}(a|x)} \cdot (R(x,a) - R^{\star}(x,a))$$

$$= -\frac{1}{\eta} \log Z_R(x) + \mathbb{E}_{\pi_R^{\eta}} \left[R(x,a) - R^{\star}(x,a) \right]. \tag{B.4}$$

Then, since $\pi_t = \pi_{\widehat{R}}^{\eta}$, the right-hand side of Equation B.3 can be written as:

$$\frac{1}{\eta} \log Z_{R^*}(x) - \frac{1}{\eta} \log Z_{\hat{R}_t}(x) + \mathbb{E}_{\pi_t} \left[\hat{R}_t(x, a) - R^*(x, a) \right] = f(x, \hat{R}_t) - f(x, R^*).$$

First, we present the lemma that gives the derivatives of π_R^{η} and Z_R , with the proof given in Appendix B.3.1.

Lemma B.4. For any $(x, a) \in \mathcal{X} \times \mathcal{A}$, we have

$$\frac{\partial \pi_R^{\eta}(x, a')}{\partial R(x, a)} = \begin{cases} \eta \pi_R^{\eta}(x, a) - \eta \pi_R^{\eta}(x, a)^2, & \text{if } a = a', \\ -\eta \pi_R^{\eta}(x, a') \pi_R^{\eta}(x, a), & \text{if } a \neq a'. \end{cases}$$

$$\frac{\partial Z_R(x)}{\partial R(x, a)} = \eta \pi_{\text{ref}}(a|x) \exp(\eta R(x, a)),$$

$$\frac{\partial \mu_R(x)}{\partial R(x, a)} = \eta \pi_R^{\eta}(x, a) \left(R(x, a) - R^{\star}(x, a) - \mu_R(x)\right) + \pi_R^{\eta}(x, a),$$

where $\mu_R(x) := \mathbb{E}_{a \sim \pi_{\mathcal{D}}^{\eta}(\cdot|x)} [R(x,a) - R^{\star}(x,a)].$

Then, we compute the derivative of f(x, R) as follows:

$$\begin{split} \frac{\partial f(x,R)}{\partial R(x,a)} &= -\frac{1}{\eta} \frac{\partial}{\partial R(x,a)} \log Z_R(x) + \frac{\partial}{\partial R(x,a)} \mathbb{E}_{\pi_R^{\eta}} \left[R(x,a) - R^{\star}(x,a) \right] \\ &= -\frac{1}{\eta} \frac{1}{Z_R(x)} \frac{\partial Z_R(x)}{\partial R(x,a)} + \frac{\partial}{\partial R(x,a)} \left[\pi_R^{\eta}(a|x) \cdot \left(R(x,a) - R^{\star}(x,a) \right) \right] \\ &+ \frac{\partial}{\partial R(x,a)} \left[\sum_{a' \neq a} \pi_R^{\eta}(a'|x) \cdot \left(R(x,a') - R^{\star}(x,a') \right) \right] \\ &= -\pi_R^{\eta}(x,a) + \pi_R^{\eta}(x,a) + \frac{\partial \pi_R^{\eta}(a|x)}{\partial R(x,a)} \cdot \left(R(x,a) - R^{\star}(x,a) \right) \\ &+ \sum_{a' \neq a} \frac{\partial \pi_R^{\eta}(a'|x)}{\partial R(x,a)} \cdot \left(R(x,a') - R^{\star}(x,a') \right) \end{aligned} \tag{Lemma B.4} \\ &= \eta \pi_R^{\eta}(a|x) \cdot \left(R(x,a) - R^{\star}(x,a) - \mathbb{E}_{a'' \sim \pi_R^{\eta}(\cdot|x)} \left[R(x,a'') - R^{\star}(x,a'') \right] \right) \\ &= \eta \pi_R^{\eta}(a|x) \cdot \left(R(x,a) - R^{\star}(x,a) - \mu_R(x) \right), \end{split} \tag{Lemma B.4}$$

where $\mu_R(x) := \mathbb{E}_{a'' \sim \pi_R^{\eta}(\cdot|x)} \left[R(x, a'') - R^{\star}(x, a'') \right]$. Note that when $R = R^{\star}$, we have $\mu_{R^{\star}}(x) = 0$, which implies

$$\frac{\partial f(x, R^{\star})}{\partial R(x, a)} = 0.$$

Moreover, the second-order gradient of f can be expressed as:

$$\frac{\partial^{2} f(x,R)}{\partial R(x,a')\partial R(x,a)}$$

$$= \frac{\partial}{\partial R(x,a')} \left(\eta \pi_{R}^{\eta}(a|x) \cdot \left(R(x,a) - R^{\star}(x,a) - \mu_{R}(x) \right) \right)$$

$$= \eta \frac{\partial \pi_{R}^{\eta}(a|x)}{\partial R(x,a')} \cdot \left(R(x,a) - R^{\star}(x,a) - \mu_{R}(x) \right) + \eta \pi_{R}^{\eta}(a|x) \cdot \left(\mathbf{1}_{a=a'} - \frac{\partial \mu_{R}(x)}{\partial R(x,a')} \right)$$

$$= \eta^{2} \pi_{R}^{\eta}(a|x) \left(\mathbf{1}_{a=a'} - \pi_{R}^{\eta}(a'|x) \right) \left(R(x,a) - R^{\star}(x,a) - \mu_{R}(x) \right)$$

$$+ \eta \pi_{R}^{\eta}(a|x) \left(\mathbf{1}_{a=a'} - \eta \pi_{R}^{\eta}(x,a') \left(R(x,a') - R^{\star}(x,a') - \mu_{R}(x) \right) + \pi_{R}^{\eta}(x,a') \right)$$
(Lemma B.4)
$$= \eta \pi_{R}^{\eta}(a|x) \left(\mathbf{1}_{a=a'} - \pi_{R}^{\eta}(a'|x) \right)$$

$$+ \eta^{2} \pi_{R}^{\eta}(a|x) \left(\mathbf{1}_{a=a'} - \pi_{R}^{\eta}(a'|x) \right) \left(R(x,a) - R^{\star}(x,a) - \mu_{R}(x) \right)$$

$$- \pi_{R}^{\eta}(a'|x) \left(R(x,a') - R^{\star}(x,a') - \mu_{R}(x) \right) \right].$$

For simplicity let $\Delta R_t = \hat{R}_t - R^\star$ and $v_t^\alpha(x,a) = \alpha \Delta R_t(x,a) - \mu_{R^\star + \alpha \Delta R_t}(x) = \alpha \Delta R_t(x,a) - \alpha \mathbb{E}_{\pi_{R^\star + \alpha \Delta R_t}^\eta}[\Delta R_t(x,a'')]$. Then, using the exact second-order Taylor expansion, we have

$$f(x, \hat{R}_{t}) - f(x, R^{*}) = f(x, R^{*} + \alpha \Delta R_{t}) - f(x, R^{*})$$

$$= \int_{0}^{1} (1 - \alpha) \left[\sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} \Delta R_{t}(x, a) \frac{\partial^{2} f(x, R^{*} + \alpha \Delta R_{t})}{\partial R(x, a') \partial R(x, a)} \Delta R_{t}(x, a') \right] d\alpha \qquad \left(\frac{\partial f(x, R^{*})}{\partial R(x, a)} = 0 \right)$$

$$= \int_{0}^{1} (1 - \alpha) \left[\eta \sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) \left(\Delta R_{t}(x, a) \right)^{2} - \eta \left(\sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) \Delta R_{t}(x, a) \right)^{2} + \eta^{2} \sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) v_{t}^{\alpha}(x, a) \left(\Delta R_{t}(x, a) \right)^{2} - 2\eta^{2} \left(\sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) v_{t}^{\alpha}(x, a) \Delta R_{t}(x, a) \right) \left(\sum_{a' \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a'|x) \Delta R_{t}(x, a') \right) \right] d\alpha.$$
(B.5)

Plugging $v_t^{\alpha}(x,a) = \alpha \Delta R_t(x,a) - \alpha \mathbb{E}_{\pi_{R^*+\alpha\Delta R_t}^{\eta}}[\Delta R_t(x,a'')]$ into the right-hand side, we can further simplify the second and third terms as follows:

$$\eta^{2} \sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) v_{t}^{\alpha}(x, a) \left(\Delta R_{t}(x, a)\right)^{2}$$

$$-2\eta^{2} \left(\sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) v_{t}^{\alpha}(x, a) \Delta R_{t}(x, a)\right) \left(\sum_{a' \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a'|x) \Delta R_{t}(x, a')\right)$$

$$= \eta^{2} \alpha \left[\sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) \left(\Delta R_{t}(x, a)\right)^{3} - 3\mathbb{E}_{\pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x, a'')\right] \sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) \left(\Delta R_{t}(x, a)\right)^{2} + 2\left(\mathbb{E}_{\pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x, a'')\right]\right)^{3}\right]$$

$$(\mathbb{E}[(X - \mathbb{E}[X])X] = \mathbb{E}[X^{2}] - (\mathbb{E}[X])^{2})$$

$$= \eta^{2} \alpha \sum_{a \in \mathcal{A}} \pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}(a|x) \left(\Delta R_{t}(x, a) - \mathbb{E}_{\pi_{R^{*} + \alpha \Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x, a'')\right]\right)^{3}.$$

$$(\mathbb{E}[(X - \mathbb{E}[X])^{3}] = \mathbb{E}[X^{3}] - 3\mathbb{E}[X]\mathbb{E}[X^{2}] + 2(\mathbb{E}[X])^{3})$$

Using this, we can rewrite the right-hand side of Equation B.5 as follows:

$$f(x, \hat{R}_t) - f(x, R^*) = \int_0^1 (1 - \alpha) \left[\eta \operatorname{Var}_t^{\alpha}(x) + \eta^2 \alpha M_t^{\alpha}(x) \right] d\alpha, \tag{B.6}$$

where we define

$$\operatorname{Var}_{t}^{\alpha}(x) := \sum_{a \in \mathcal{A}} \pi_{R^{\star} + \alpha \Delta R_{t}}^{\eta}(a|x) \left(\Delta R_{t}(x, a)\right)^{2} - \left(\sum_{a \in \mathcal{A}} \pi_{R^{\star} + \alpha \Delta R_{t}}^{\eta}(a|x) \Delta R_{t}(x, a)\right)^{2}$$
$$M_{t}^{\alpha}(x) := \sum_{a \in \mathcal{A}} \pi_{R^{\star} + \alpha \Delta R_{t}}^{\eta}(a|x) \left(\Delta R_{t}(x, a) - \mathbb{E}_{\pi_{R^{\star} + \alpha \Delta R_{t}}^{\eta}}[\Delta R_{t}(x, a'')]\right)^{3}.$$

The following lemma is a useful tool for calculating the right-hand side of Equation B.6. Its proof is presented in Appendix B.3.2.

Lemma B.5. Let $\pi_{\alpha}(a|x) := \frac{\pi_{\text{ref}}(a|x) \exp(\eta R_{\alpha}(x,a))}{Z_{\alpha}(x)}$, where $R_{\alpha} = R^{\star} + \alpha \Delta R$ with $R^{\star}, \Delta R \in \mathbb{R}$, and $Z_{\alpha}(x) = \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|x) \exp(\eta R_{\alpha}(x,a))$. Then, for any $(x,a) \in \mathcal{X} \times \mathcal{A}$, we have

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \pi_{\alpha}(a|x) = \eta \pi_{\alpha}(a|x) \left(\Delta R(x,a) - \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)] \right),$$

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)] = \eta \mathbb{E}_{\pi_{\alpha}} \left[\left(\Delta R(x,a) - \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)] \right)^{2} \right],$$

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)^{2}] = \eta \left(\mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)^{3}] - \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)^{2}] \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)] \right).$$

Then, by Lemma B.5, we show that

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \operatorname{Var}_{t}^{\alpha}(x) \\
= \frac{\mathrm{d}}{\mathrm{d}\alpha} \left(\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[(\Delta R_{t}(x,a))^{2} \right] - \left(\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x,a) \right] \right)^{2} \right) \\
= \frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[(\Delta R_{t}(x,a))^{2} \right] - 2\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x,a) \right] \cdot \frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x,a) \right] \\
= \eta \left(\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[(\Delta R_{t}(x,a))^{3} \right] - \mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x,a) \right] \mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[(\Delta R_{t}(x,a))^{2} \right] \\
- 2\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x,a) \right] \cdot \operatorname{Var}_{t}^{\alpha}(x) \right) \qquad (\text{Lemma B.5}) \\
= \eta \left(\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[(\Delta R_{t}(x,a))^{3} \right] - 3\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x,a) \right] \mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[(\Delta R_{t}(x,a))^{2} \right] \\
+ 2 \left(\mathbb{E}_{\pi_{R^{*}+\alpha\Delta R_{t}}^{\eta}} \left[\Delta R_{t}(x,a) \right] \right)^{3} \right) \qquad (\text{Definition of Var}_{t}^{\alpha}(x)) \\
= \eta M_{t}^{\alpha}(x). \qquad (\text{Definition of } M_{t}^{\alpha}(x)) \right)$$

Therefore, Equation B.6 can be further simplified as:

1171
1172
1173
1174
1175
$$= \eta \left[\int_0^1 (1-\alpha) \left[\eta \operatorname{Var}_t^{\alpha}(x) + \eta^2 \alpha M_t^{\alpha}(x) \right] d\alpha \right]$$
1174
1175
$$= \eta \left[\int_0^1 (1-\alpha) \operatorname{Var}_t^{\alpha}(x) d\alpha + \int_0^1 \alpha \frac{\mathrm{d}}{\mathrm{d}\alpha} \operatorname{Var}_t^{\alpha}(x) d\alpha \right]$$
1176
1177
$$= \eta \left[\int_0^1 (1-\alpha) \operatorname{Var}_t^{\alpha}(x) d\alpha + \left[\alpha \operatorname{Var}_t^{\alpha}(x) \right]_0^1 - \int_0^1 \operatorname{Var}_t^{\alpha}(x) d\alpha \right]$$
(integration by parts)
1180
$$= \eta \left[\operatorname{Var}_t^{\alpha=1}(x) - \int_0^1 \alpha \operatorname{Var}_t^{\alpha}(x) d\alpha \right]$$
1182
$$= \eta \mathbb{E}_{\pi_{\widehat{R}_t}^{\eta}} \left[\left(\Delta R_t(x, a) - \mathbb{E}_{\pi_{\widehat{R}_t}^{\eta}} [\Delta R_t(x, a)] \right)^2 \right] - \eta \int_0^1 \alpha \operatorname{Var}_t^{\alpha}(x) d\alpha$$
1184
$$\leq \eta \mathbb{E}_{\pi_{\widehat{R}_t}^{\eta}} \left[\left(\Delta R_t(x, a) - \mathbb{E}_{\pi_{\widehat{R}_t}^{\eta}} [\Delta R_t(x, a)] \right)^2 \right] \qquad (\operatorname{Var}_t^{\alpha}(x) \geq 0)$$

$$\leq \eta \mathbb{E}_{\pi_{\widehat{R}_t}^{\eta}} \left[\left(\Delta R_t(x, a) - \mathbb{E}_{\pi_{\widehat{R}_t}^{\eta}} [\Delta R_t(x, a)] \right)^2 \right] \leq \mathbb{E}[X^2]$$
1187

Recall that $\pi_t = \pi_{\hat{R}_t}^{\eta}$ and $\Delta R_t = \hat{R}_t - R^{\star}$. Hence, we obtain

1190
1191
$$J_t^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_t^{\eta}(\pi_t, R^{\star}) \leq \eta \mathbb{E}_{a \sim \pi_t(\cdot | x_t)} \left[\left(\hat{R}_t(x_t, a) - R^{\star}(x_t, a) \right)^2 \right].$$
1192

This concludes the proof of Lemma B.2.

B.2.2 PROOF OF LEMMA B.3

Proof of Lemma B.3. For simple presentation, we write $\mathbb{E}_{\pi}[\cdot] = \mathbb{E}_{a \sim \pi(\cdot \mid x)}[\cdot]$. Then, for any $t \in [T]$, we have

$$\mathbb{E}_{a \sim \pi^{\star}(\cdot|x_{t})}[R^{\star}(x_{t}, a)] = J_{t}^{\eta}(\pi^{\star}, R^{\star}) + \frac{1}{\eta} \operatorname{KL}(\pi^{\star}(\cdot|x_{t})|\pi_{\operatorname{ref}}(\cdot|x_{t})) \qquad \text{(Definition of } J_{t}^{\eta})$$

$$\leq J_{t}^{\eta}(\pi_{\eta}^{\star}, R^{\star}) + \frac{1}{\eta} \operatorname{KL}(\pi^{\star}(\cdot|x_{t})|\pi_{\operatorname{ref}}(\cdot|x_{t})). \qquad \text{(Definition of } \pi_{\eta}^{\star})$$

Moreover, since the KL divergence is always non-negative, we get

$$\mathbb{E}_{a \sim \pi_t(\cdot | x_t)} [R^{\star}(x_t, a)] \geqslant \mathbb{E}_{a \sim \pi_t(\cdot | x_t)} [R^{\star}(x_t, a)] - \frac{1}{\eta} \operatorname{KL} (\pi_t(\cdot | x_t) | \pi_{\operatorname{ref}}(\cdot | x_t))$$
$$= J_t^{\eta}(\pi_t, R^{\star}).$$

Combining the above two results, we obtain

$$\mathbb{E}_{a \sim \pi^{\star}(\cdot|x_{t})}[R^{\star}(x_{t}, a)] - \mathbb{E}_{a \sim \pi_{t}(\cdot|x_{t})}[R^{\star}(x_{t}, a)]$$

$$\leq J_{t}^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_{t}^{\eta}(\pi_{t}, R^{\star}) + \frac{1}{\eta} \operatorname{KL}(\pi^{\star}(\cdot|x_{t})||\pi_{\operatorname{ref}}(\cdot|x_{t})),$$

which concludes the proof of Lemma B.3.

B.3 Supporting Results for Lemma B.2

B.3.1 PROOF OF LEMMA B.4

Proof of Lemma B.4. First, we compute the derivative of the policy $\pi_R^{\eta}(x, a)$. For any $(x, a) \in \mathcal{X} \times \mathcal{A}$, we have

$$\begin{split} \frac{\partial \pi_R^{\eta}(x,a)}{\partial R(x,a)} &= \frac{\partial}{\partial R(x,a)} \left(\frac{1}{Z_R(x)} \pi_{\text{ref}}(a|x) \exp(\eta R(x,a)) \right) \\ &= \frac{\eta \pi_{\text{ref}}(a|x) \exp(\eta R(x,a))}{Z_R(x)} - \frac{\pi_{\text{ref}}(a|x) \exp(\eta R(x,a))}{Z_R(x)^2} \cdot \frac{\partial Z_R(x)}{\partial R(x,a)} \\ &= \frac{\eta \pi_{\text{ref}}(a|x) \exp(\eta R(x,a))}{Z_R(x)} - \frac{\pi_{\text{ref}}(a|x) \exp(\eta R(x,a))}{Z_R(x)^2} \cdot \eta \pi_{\text{ref}}(a|x) \exp(\eta R(x,a)) \\ &= \eta \pi_R^{\eta}(x,a) - \eta \pi_R^{\eta}(x,a)^2. \end{split}$$

Moreover, for any $(x, a, a') \in \mathcal{X} \times \mathcal{A} \times \mathcal{A}$ with $a' \neq a$, we obtain

$$\begin{split} \frac{\partial \pi_R^{\eta}(x,a')}{\partial R(x,a)} &= \pi_{\text{ref}}(a'|x) \exp \left(\eta R(x,a') \right) \cdot \frac{\partial}{\partial R(x,a)} \left(\frac{1}{Z_R(x)} \right) \\ &= -\frac{\pi_{\text{ref}}(a'|x) \exp \left(\eta R(x,a') \right)}{Z_R(x)^2} \cdot \eta \pi_{\text{ref}}(a|x) \exp \left(\eta R(x,a) \right) \\ &= -\eta \pi_R^{\eta}(x,a') \pi_R^{\eta}(x,a). \end{split}$$

Next, we compute the derivative of $Z_R(x)$. For any $(x, a) \in \mathcal{X} \times \mathcal{A}$, we get

$$\frac{\partial Z_R(x)}{\partial R(x,a)} = \frac{\partial}{\partial R(x,a)} \left(\mathbb{E}_{\pi_{\text{ref}}} [\exp(\eta R(x,a))] \right) = \eta \pi_{\text{ref}}(a|x) \exp(\eta R(x,a)),$$

Finally, we compute the derivative of $\mu_R(x) = \mathbb{E}_{a \sim \pi_R^{\eta}(\cdot|x)} [R(x,a) - R^{\star}(x,a)]$. For any $(x,a) \in \mathcal{X} \times \mathcal{A}$, we have

$$\frac{\partial \mu_R(x)}{\partial R(x,a)} = \sum_{a' \in \mathcal{A}} \frac{\partial \pi_R^{\eta}(a'|x)}{\partial R(x,a)} \left(R(x,a') - R^{\star}(x,a') \right) + \pi_R^{\eta}(a|x)
= \eta \pi_R^{\eta}(a|x) \sum_{a' \in \mathcal{A}} \left(\mathbf{1}_{a=a'} - \pi_R^{\eta}(a'|x) \right) \cdot \left(R(x,a') - R^{\star}(x,a') \right) + \pi_R^{\eta}(a|x)
= \eta \pi_R^{\eta}(x,a) \left(R(x,a) - R^{\star}(x,a) - \mu_R(x) \right) + \pi_R^{\eta}(x,a).$$

Thus, we conclude the proof of Lemma B.4.

B.3.2 PROOF OF LEMMA B.5

Proof of Lemma B.5. For the first property, a simple calculation gives

$$\frac{\mathrm{d}}{\mathrm{d}\alpha}\pi_{\alpha}(a|x) = \frac{\pi_{\mathrm{ref}}(a|x)\exp\left(\eta R_{\alpha}(x,a)\right) \cdot \eta \Delta R(x,a) Z_{\alpha}(x) - \pi_{\mathrm{ref}}(a|x)\exp\left(\eta R_{\alpha}(x,a)\right) \cdot \frac{\mathrm{d}Z_{\alpha}(x)}{\mathrm{d}\alpha}}{Z_{\alpha}(x)^{2}}$$

$$= \frac{\pi_{\mathrm{ref}}(a|x)\exp\left(\eta R_{\alpha}(x,a)\right)}{Z_{\alpha}(x)} \left[\eta \Delta R(x,a) - \frac{1}{Z_{\alpha}(x)} \frac{\mathrm{d}Z_{\alpha}(x)}{\mathrm{d}\alpha}\right]$$

$$= \pi_{\alpha}(a|x) \left[\eta \Delta R(x,a) - \frac{1}{Z_{\alpha}(x)} \frac{\mathrm{d}Z_{\alpha}(x)}{\mathrm{d}\alpha}\right].$$
(B.7)

Moreover, we get

$$\frac{\mathrm{d}Z_{\alpha}(x)}{\mathrm{d}\alpha} = \sum_{a \in \mathcal{A}} \pi_{\mathrm{ref}}(a|x) \exp\left(\eta R_{\alpha}(x,a)\right) \cdot \eta \Delta R(x,a) = \eta Z_{\alpha}(x) \sum_{a \in \mathcal{A}} \pi_{\alpha}(a|x) \Delta R(x,a)$$

$$= \eta Z_{\alpha}(x) \mathbb{E}_{\pi_{\alpha}}[\Delta R(x,a)]. \tag{B.8}$$

Plugging Equation B.8 into Equation B.7, we obtain the first property.

Now, we prove the second property.

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)] &= \sum_{a \in \mathcal{A}} \frac{\mathrm{d}\pi_{\alpha}(a|x)}{\mathrm{d}\alpha} \Delta R(x,a) \\ &= \eta \sum_{a \in \mathcal{A}} \pi_{\alpha}(a|x) \big(\Delta R(x,a) - \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)] \big) \Delta R(x,a) \quad \text{(first property)} \\ &= \eta \left(\mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)^2] - (\mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)])^2 \right) \\ &= \eta \mathbb{E}_{\pi_{\alpha}} \left[(\Delta R(x,a) - \mathbb{E}_{\pi_{\alpha}} [\Delta R(x,a)])^2 \right]. \end{split}$$

Similarly, substituting $\Delta R(x, a)$ with $\Delta R(x, a)^2$ in the above analysis, we obtain

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbb{E}_{\pi_{\alpha}} [\Delta R(x, a)^{2}] = \eta \sum_{a \in \mathcal{A}} \pi_{\alpha}(a|x) (\Delta R(x, a) - \mathbb{E}_{\pi_{\alpha}} [\Delta R(x, a)]) \Delta R(x, a)^{2} \quad \text{(first property)}$$

$$= \eta \left(\mathbb{E}_{\pi_{\alpha}} [\Delta R(x, a)^{3}] - \mathbb{E}_{\pi_{\alpha}} [\Delta R(x, a)^{2}] \mathbb{E}_{\pi_{\alpha}} [\Delta R(x, a)] \right),$$

which proves the last property.

B.4 DISCUSSION ON SPECIFIC FUNCTION CLASSES

In this subsection, we supplement the result of Theorem 1 by providing a more detailed discussion of the tightness of our (unregularized) regret bound for several special function classes. We set the reference policy to be uniform, i.e., $\pi_{\rm ref} = {\rm Unif}(\mathcal{A})$. Then, for any policy π , it holds that ${\rm KL}(\pi \| \pi_{\rm ref}) = \sum_a \left(\pi(a) \log \pi(a) - \pi(a) \log \frac{1}{|\mathcal{A}|}\right) \leqslant \log |\mathcal{A}| = \log N$. Hence, KL-EXP yields the following regret bounds for special function classes:

- **1. Linear classes:** When $R^{\star} \in \mathcal{R}$ and the reward function class \mathcal{R} is linear, i.e., $\mathcal{R} = \{R: R = \phi(x, a)^{\top}\theta, \theta \in \mathbb{R}^d, \|\theta\|_2 \leqslant 1\}$, where $\phi(x, a) \in \mathbb{R}^d$ is a known feature map satisfying $\|\phi(x, a)\|_2 \leqslant 1$, the Vovk-Azoury-Warmuth forecaster (Vovk, 1997; Azoury & Warmuth, 2001) guarantees $\operatorname{Reg}_{\operatorname{Sq}}(T) = \mathcal{O}(d\log(T/d))$ (Example 1), which implies $\operatorname{Regret}(T) = \mathcal{O}(\sqrt{dT\log N\log T})$. As stated in Remark 3, this bound is minimax-optimal, matching the lower bound $\Omega(\sqrt{dT\log N\log(T/d)})$ (Li et al., 2019) up to logarithmic d factors. It is remarkable that we obtain this $\tilde{\mathcal{O}}(\sqrt{dT\log N})$ -type regret bound without relying on the difficult-to-implement "layered data partitioning" technique required in prior works (Auer, 2002; Chu et al., 2011; Li et al., 2019). Our algorithm is simple to implement: it only requires solving the KL-regularized objective in Equation 1 (with the closed-form solution in Equation 2) using the reward estimator \hat{R}_t returned by the online regression oracle. We believe this opens a promising direction for developing algorithms that are both practical and statistically optimal in linear contextual bandits.
- **2. Multi-armed bandits (MABs):** The MAB case follows directly from the linear case by setting d = N. We achieve $\text{Reg}_{So}(T) = \mathcal{O}(N \log(T/N))$ and $\text{Regret}(T) = \mathcal{O}(\sqrt{NT \log N \log(T/N)})$.
- 3. Generalized linear models (GLMs): For GLM reward function class, i.e., $\mathcal{R} = \{R : R = \mu(\phi(x,a)^\top\theta), \theta \in \mathbb{R}^d, \|\theta\|_2 \le 1\}$, where $\mu : \mathbb{R} \to [0,1]$ is a fixed non-decreasing 1-Lipschitz link function and $\phi(x,a) \in \mathbb{R}^d$ is a known feature map with $\|\phi(x,a)\|_2 \le 1$, if $R^\star \in \mathcal{R}$, the GLMtron algorithm (Kakade et al., 2011) guarantees $\operatorname{Reg}_{\operatorname{Sq}}(T) = \mathcal{O}(\kappa_\mu^2 d \log(T/d))$, where $1/\dot{\mu} \le \kappa_\mu$. This, in turn, implies $\operatorname{Regret}(T) = \mathcal{O}\left(\kappa_\mu \sqrt{dT \log N \log T}\right)$, which is tighter than the bound $\mathcal{O}\left(\kappa_\mu(\log T)^{1.5}\sqrt{dT \log N}\right)$ (Li et al., 2017) by a factor of $\log T$. On the other hand, Lee et al. (2024); Sawarni et al. (2024) establish a κ_μ -improved regret bound of $\tilde{\mathcal{O}}\left(d\sqrt{T/\kappa_\mu^\star}\right)$, where $\kappa_\mu^\star := \frac{1}{\dot{\mu}((x^\star)^\top\theta^\star)}$, though with a looser dependence on \sqrt{d} than ours. It remains an open question whether a $\tilde{\mathcal{O}}(\sqrt{dT \log N})$ -type regret bound can be attained while simultaneously improving the dependence on κ_μ .
- **4. Bounded eluder dimension:** Under the realizability assumption (Assumption 1), i.e., $R^* \in \mathcal{R}$, and the reward function class \mathcal{R} has bounded eluder dimension (Definition C.1), the empirical risk minimization (ERM) algorithm achieves, with probability at least 1δ , $\mathrm{Reg}_{\mathrm{Sq}}(T) = \mathcal{O}\big(d_{\mathrm{E}}\log(\mathcal{N}_{\mathcal{R}}(\epsilon)T)\big)$ (Lemma C.2). Consequently, we obtain the unregularized regret bound $\mathrm{Regret}(T) = \mathcal{O}\big(\sqrt{d_{\mathrm{E}}T\log N\log(\mathcal{N}_{\mathcal{R}}(\epsilon)T)}\big)$. In comparison, the existing bound of Russo & Van Roy (2013) is $\mathcal{O}\big(\sqrt{d_{\mathrm{E}}T\log(\mathcal{N}_{\mathcal{R}}(\epsilon)T)}\big)$, which shows that our result is tight up to a $\sqrt{\log N}$ factor.

C CASE: \mathcal{R} WITH BOUNDED ELUDER DIMENSION (REMARK 1)

In this subsection, we analyze the setting where the reward function class \mathcal{R} has bounded eluder dimension (Russo & Van Roy, 2013), in order to enable a direct comparison with prior work (Zhao et al., 2025).

We define the uncertainty and eluder dimension, following Zhao et al. (2025).

Definition C.1. For any sequence $\mathcal{D}_t = \{(x_s, a_s)\}_{s=1}^{t-1}$, we define the uncertainty of (x, a) with respect to \mathcal{R} as:

$$U_{\mathcal{R},\lambda}(x,a;\mathcal{D}_t) := \sup_{R_1,R_2 \in \mathcal{R}} \frac{|R_1(x,a) - R_2(x,a)|}{\sqrt{\lambda + \sum_{s=1}^{t-1} (R_1(x_s,a_s) - R_2(x_s,a_s))^2}}.$$

And the eluder dimension is defined as:

$$d_{\mathcal{E}} := \sup_{x_{1:T}, a_{1:T}} \sum_{t=1}^{T} \min \left\{ 1, U_{\mathcal{R}, \lambda}(x_t, a_t; \mathcal{D}_t)^2 \right\}. \tag{C.1}$$

We also define the confidence set \mathcal{R}_t as follows:

$$\mathcal{R}_t := \left\{ R \in \mathcal{R} : \sum_{s=1}^{t-1} \left(R(x_s, a_s) - \hat{R}_t(x_s, a_s) \right)^2 + \lambda \leqslant \beta_T^2 = 16 \log(\mathcal{N}_{\mathcal{R}}(\epsilon) T / \delta) \right\},\,$$

where $\lambda > 0$. We can then bound the estimation error using the following lemma.

Lemma C.1 (Lemma 4.5 of Zhao et al. 2025). Let \hat{R}_t be the empirical risk minimizer (ERM), i.e., $\hat{R}_t \leftarrow \underset{s=1}{\operatorname{argmin}}_{R \in \mathcal{R}} \sum_{s=1}^{t-1} (R(x_s, a_s) - y_s)^2$. Then, under Assumption 1 and the condition that the noises ϵ_t are conditional 1-subGaussian, we have with probability at least $1 - \delta$, for all $t \in [T]$, we have

$$\hat{R}_t(x, a) - R^*(x, a) \leq \min\{1, \beta_T \cdot U_{\mathcal{R}_t, \lambda}(x, a; \mathcal{D}_t)\}, \quad \forall (x, a) \in \mathcal{X} \times \mathcal{A}.$$

The following lemma is useful for the subsequent analysis.

Lemma C.2. Under Assumption 1, if OracleSq is chosen as the standard ERM algorithm, then with probability at least $1 - \delta$ we obtain

$$\sum_{t=1}^{T} \left(\hat{R}_t(x_t, a_t) - r_t \right)^2 - \sum_{t=1}^{T} \left(R^{\star}(x_t, a_t) - r_t \right)^2 = \mathcal{O}\left(d_{\mathsf{E}} \log(\mathcal{N}_{\mathcal{R}}(\epsilon)T) \right).$$

Proof of Lemma C.2. Let $M_t:=(\hat{R}_t(x_t,a_t)-r_t)^2-(R^\star(x_t,a_t)-r_t)^2$ and $Z_t:=M_t-\mathbb{E}[M_t\mid\mathcal{F}_{t-1}]$. We define the filtration $\mathcal{F}_{t-1}=\sigma(x_1,a_1,r_1,\ldots,x_{t-1},a_{t-1},r_{t-1},x_t)$. Then, by Lemma B.1 and Freedman's inequality (Lemma G.1) with $\beta=1/8$, with probability at least $1-\delta$, we have

$$\begin{split} \sum_{t=1}^{T} M_{t} & \leqslant \sum_{t=1}^{T} \mathbb{E}[M_{t} | \mathcal{F}_{t=1}] + \frac{1}{8} \sum_{t=1}^{T} \mathbb{E}[Z_{t}^{2} | \mathcal{F}_{t-1}] + 8 \log \frac{1}{\delta} & \text{(Lemma G.1, w.p. } 1 - \delta) \\ & \leqslant \frac{3}{2} \sum_{t=1}^{T} \mathbb{E}[M_{t} | \mathcal{F}_{t-1}] + 8 \log \frac{1}{\delta} & \text{(Lemma B.1)} \\ & = \frac{3}{2} \sum_{t=1}^{T} \mathbb{E}_{a \sim \pi_{t}} \left[(\hat{R}_{t}(x_{t}, a_{t}) - R^{\star}(x_{t}, a_{t}))^{2} | \mathcal{F}_{t-1} \right] + 8 \log \frac{1}{\delta} \\ & \leqslant 3 \sum_{t=1}^{T} \left(\hat{R}_{t}(x_{t}, a_{t}) - R^{\star}(x_{t}, a_{t}) \right)^{2} + 16 \log \frac{2}{\delta}. & \text{(Lemma G.2, w.p. } 1 - \delta) \end{split}$$

Hence, we derive

$$\begin{split} &\sum_{t=1}^T \left(\hat{R}_t(x_t, a_t) - r_t\right)^2 - \sum_{t=1}^T \left(R^\star(x_t, a_t) - r_t\right)^2 \\ &\leqslant 3 \sum_{t=1}^T \left(\hat{R}_t(x_t, a_t) - R^\star(x_t, a_t)\right)^2 + 16\log\frac{2}{\delta} \\ &\leqslant 3\beta_T^2 \sum_{t=1}^T \min\left\{1, U_{\mathcal{R}_t, \lambda}(x_t, a_t; \mathcal{D}_t)^2\right\} + 16\log\frac{2}{\delta} \\ &\leqslant 48d_{\mathrm{E}} \log(\mathcal{N}_{\mathcal{R}}(\epsilon)T/\delta) + 16\log\frac{2}{\delta}. \end{split} \tag{Lemma C.1, w.p. } 1 - \delta)$$

By setting $\delta \leftarrow \frac{\delta}{3}$, the proof is complete.

We now present the claim in Remark 1 more formally.

Proposition C.1 (Regret under bounded eluder dimension). Suppose the eluder dimension defined in Equation C.1 is finite. Let the online regression oracle OracleSq be the ERM predictor. Under Assumptions 1 and 2, for any $\delta > 0$, KL-EXP (Algorithm 1) guarantees that with probability at least $1 - \delta$,

$$\mathbf{Regret}_{\mathrm{KL}}(T,\eta) = \mathcal{O}\left(\eta d_{\mathrm{E}}\log\left(\mathcal{N}_{\mathcal{R}}(\epsilon)T\right)\right), \quad and \quad \mathbf{Regret}(T) = \mathcal{O}\left(\eta d_{\mathrm{E}}\log\left(\mathcal{N}_{\mathcal{R}}(\epsilon)T\right) + \frac{DT}{\eta}\right),$$

$$where \ D := \frac{1}{T} \sum_{t=1}^{T} \mathrm{KL}\left(\pi^{\star}(\cdot \| x_{t}) \| \pi_{\mathrm{ref}}(\cdot \| x_{t})\right).$$

Algorithm D.1 0EP0 (Oracle-Efficient Policy Optimization)

- 1: **Inputs:** regularization parameter η , reference policy π_{ref} , online regression oracle OracleLog.
- 2: **Initialize:** choose any $R_1 \in \mathcal{R}$.
- 3: **for** round t = 1 to T **do**
- 4: Observe context $x_t \in \mathcal{X}$.
 - 5: Compute policy $\pi_t(\cdot|x_t) \propto \pi_{\text{ref}}(\cdot|x_t) \exp(\eta \hat{R}_t(x_t,\cdot))$ via Equation 2.
- 6: Sample action $a_t^1, a_2^t \sim \pi_t(\cdot|x_t)$ and receive preference feedback y_t .
 - 7: Update \hat{R}_{t+1} for the next round using OracleLog via Equation 7.
 - 8: end for

Proof of Proposition C.1. Then, following a similar analysis to the proof of Theorem 1, we can bound the regret as follows:

$$\begin{aligned} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) &= \sum_{t=1}^{T} J_{t}^{\eta}(\pi_{\eta}^{\star},R^{\star}) - J_{t}^{\eta}(\pi_{t},R^{\star}) \\ &\leqslant \eta \sum_{t=1}^{T} \mathbb{E}_{a_{t} \sim \pi_{t}(\cdot|x_{t})} \left[\left(\hat{R}_{t}(x_{t},a_{t}) - R^{\star}(x_{t},a_{t}) \right)^{2} \right] \end{aligned} \qquad \text{(Lemma B.2)} \\ &\leqslant 2\eta \left[\sum_{t=1}^{T} \left(\hat{R}_{t}(x_{t},a_{t}) - r_{t} \right)^{2} - \sum_{t=1}^{T} \left(R^{\star}(x_{t},a_{t}) - r_{t} \right)^{2} \right] + 16 \log \frac{1}{\delta} \\ &\qquad \qquad \text{(Lemma G.1 and B.1 w.p. } 1 - \delta) \\ &= \mathcal{O} \big(\eta d_{\mathrm{E}} \log \left(\mathcal{N}_{\mathcal{R}}(\epsilon) T \right) \big). \end{aligned} \qquad \text{(Lemma C.2 w.p. } 1 - \delta)$$

Setting $\delta \leftarrow \frac{\delta}{2}$ yields the bound for **Regret**_{KL} (T, η) .

The bound for $\mathbf{Regret}(T)$ then follows directly from Lemma B.3. Thus, the proof of Proposition C.1 is complete.

D PROOF OF THEOREM 3

In this section, we present the proof of Theorem 3.

D.1 Main Proof of Theorem 3

We begin by introducing the key lemmas used to prove Theorem 3.

Lemma D.1. With probability at least $1 - \delta$, we have

$$\sum_{t=1}^{T} \left(\left[R^{\star}(x_{t}, a_{t}^{1}) - \hat{R}_{t}(x_{t}, a_{t}^{1}) \right] - \left[R^{\star}(x_{t}, a_{t}^{2}) - \hat{R}_{t}(x_{t}, a_{t}^{2}) \right] \right)^{2}$$

$$\leq \kappa^{2} \left(\sum_{t=1}^{T} \ell_{t}(\hat{R}_{t}) - \sum_{t=1}^{T} \ell_{t}(R^{\star}) \right) + 2\kappa^{2} \log \frac{1}{\delta}.$$

The proof is deferred to Appendix D.2.1.

Lemma D.2 (Second-order regret decomposition with baseline). For any $t \in [T]$ and any $g : \mathcal{X} \to \mathbb{R}$, we have

$$J_t^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_t^{\eta}(\pi_t, R^{\star}) \leqslant \eta \mathbb{E}_{a \sim \pi_t(\cdot | x_t)} \left[\left(\hat{R}_t(x_t, a) - R^{\star}(x_t, a) + g(x_t) \right)^2 \right].$$

The proof is deferred to Appendix D.2.2.

We now provide the proof of Theorem 3.

Proof of Theorem 3. By applying Lemma D.2 with setting

$$g_t(x) = -\mathbb{E}_{a^2 \sim \pi_t(\cdot|x)} \left[\hat{R}_t(x, a^2) - R^{\star}(x, a^2) \right]$$

we have

$$\begin{aligned} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) &= \sum_{t=1}^{T} J_{t}^{\eta}(\pi_{\eta}^{\star},R^{\star}) - J_{t}^{\eta}(\pi_{t},R^{\star}) \\ &\leqslant \eta \sum_{t=1}^{T} \mathbb{E}_{a^{1},a^{2} \sim \pi_{t}(\cdot|x_{t})} \left[\left(\hat{R}_{t}(x_{t},a^{1}) - R^{\star}(x_{t},a^{1}) - \left(\hat{R}_{t}(x_{t},a^{2}) - R^{\star}(x_{t},a^{2}) \right) \right)^{2} \right] \end{aligned} \tag{Lemma D.2}$$

$$&\leqslant 2\eta \sum_{t=1}^{T} \left(\hat{R}_{t}(x_{t},a_{t}^{1}) - R^{\star}(x_{t},a_{t}^{1}) - \left(\hat{R}_{t}(x_{t},a_{t}^{2}) - R^{\star}(x_{t},a_{t}^{2}) \right) \right)^{2} + 32\eta \log \frac{2}{\delta} \tag{Lemma G.2, w.p. } 1 - \delta)$$

$$&\leqslant 2\eta \kappa^{2} \left(\sum_{t=1}^{T} \ell_{t}(\hat{R}_{t}) - \sum_{t=1}^{T} \ell_{t}(R^{\star}) \right) + 4\eta \kappa^{2} \log \frac{1}{\delta} + 32\eta \log \frac{2}{\delta} \tag{Lemma D.1, w.p. } 1 - \delta)$$

$$&\leqslant 2\eta \kappa^{2} \mathrm{Reg}_{\mathrm{Log}}(T) + 4\eta \kappa^{2} \log \frac{1}{\delta} + 32\eta \log \frac{2}{\delta} \tag{Assumption 3}$$

(Assumption 3)

By setting $\delta \leftarrow \frac{\delta}{2}$, we establish the bound for **Regret**_{KL} (T, η) .

Furthermore, the bound on $\mathbf{Regret}(T)$ follows immediately from Lemma B.3, using the same analysis as in the proof of Theorem 1. Hence, this completes the proof of Theorem 3.

PROOFS OF LEMMAS FOR THEOREM 3

D.2.1PROOF OF LEMMA D.1

Proof of Lemma D.1. The proof of Lemma D.1 follows the analysis of Lemma D.1 in Zhao et al. (2024). However, unlike Zhao et al. (2024), where the estimator R is fixed for all t, our setting accommodates a time-varying sequence $\{\hat{R}_t\}_{t=1}^T$.

For completeness, we present the full proof below.

 For simplicity, we write $p_t^\star = \sigma \left(R^\star(x_t, a_t^1) - R^\star(x_t, a_t^2) \right)$ and $p_t = \sigma \left(\hat{R}_t(x_t, a_t^1) - \hat{R}_t(x_t, a_t^2) \right)$. We define

$$X_t := \frac{1}{2} \left(\ell_t(R^*) - \ell_t(\hat{R}_t) \right) = -\frac{1}{2} \left(y_t \log \frac{p_t^*}{p_t} + (1 - y_t) \log \frac{1 - p_t^*}{1 - p_t} \right).$$

Then, by Lemma G.3, with probability at least $1 - \delta$, we have

$$\frac{1}{2} \left(\sum_{t=1}^{T} \ell_{t}(R^{\star}) - \sum_{t=1}^{T} \ell_{t}(\hat{R}_{t}) \right) = \sum_{t=1}^{T} X_{t} \leqslant \sum_{t=1}^{T} \log \left(\mathbb{E}_{t-1}[e^{X_{t}}] \right) + \log \frac{1}{\delta}$$
 (Lemma G.3)
$$= \sum_{t=1}^{T} \log \left(p_{t}^{\star} \left(\frac{p_{t}^{\star}}{p_{t}} \right)^{-1/2} + (1 - p_{t}^{\star}) \left(\frac{1 - p_{t}^{\star}}{1 - p_{t}} \right)^{-1/2} \right) + \log \frac{1}{\delta}$$

$$= \sum_{t=1}^{T} \log \left(\sqrt{p_{t}^{\star} p_{t}} + \sqrt{(1 - p_{t}^{\star})(1 - p_{t})} \right) + \log \frac{1}{\delta}$$

$$\leqslant \sum_{t=1}^{T} \left(\sqrt{p_{t}^{\star} p_{t}} + \sqrt{(1 - p_{t}^{\star})(1 - p_{t})} - 1 \right) + \log \frac{1}{\delta}$$

$$(\log x \leqslant x - 1, \text{ for } x > 0)$$

$$= -\frac{1}{2} \sum_{t=1}^{T} \left[\left(\sqrt{p_{t}^{\star}} - \sqrt{p_{t}} \right)^{2} + \left(\sqrt{1 - p_{t}^{\star}} - \sqrt{1 - p_{t}} \right)^{2} \right] + \log \frac{1}{\delta}$$

$$(1 = \frac{1}{2} (p_{t}^{\star} + (1 - p_{t}^{\star}) + p_{t} + (1 - p_{t})))$$

$$\leqslant -\frac{1}{2} \sum_{t=1}^{T} (p_{t}^{\star} - p_{t})^{2} + \log \frac{1}{\delta}.$$
 (D.1)

where the last inequality follows from the fact that, for any $p, q \in [0, 1]$, $(\sqrt{p} - \sqrt{q})^2 + (\sqrt{1-p} - \sqrt{1-q})^2 \ge (p-q)^2$.

Now, consider the term $p_t^{\star} - p_t$. For simplicity, let $\Delta_t^{\star} = R^{\star}(x_t, a_t^1) - R^{\star}(x_t, a_t^2)$ and $\Delta_t = \hat{R}_t(x_t, a_t^1) - \hat{R}_t(x_t, a_t^2)$. Then, by the mean value theorem, we obtain

$$\begin{split} p_t^{\star} - p_t &= \sigma(\Delta_t^{\star}) - \sigma(\Delta_t) \\ &= (\Delta_t^{\star} - \Delta_t) \int_0^1 \dot{\sigma} \left(\Delta_t + \tau(\Delta_t^{\star} - \Delta_t)\right) \mathrm{d}\tau \qquad \qquad \text{(mean value theorem)} \\ &\geqslant \frac{1}{\kappa} \left(\Delta_t^{\star} - \Delta_t\right). \qquad \qquad (\dot{\sigma}(z) \geqslant \frac{1}{\kappa}, \text{Definition of } \kappa) \end{split}$$

Hence, substituting the above result into Equation D.1 and rearranging terms, we obtain

$$\sum_{t=1}^{T} \left(\left[R^{\star}(x_{t}, a_{t}^{1}) - R^{\star}(x_{t}, a_{t}^{2}) \right] - \left[\hat{R}_{t}(x_{t}, a_{t}^{1}) - \hat{R}_{t}(x_{t}, a_{t}^{2}) \right] \right)^{2}$$

$$\leq \kappa^{2} \left(\sum_{t=1}^{T} \ell_{t}(\hat{R}_{t}) - \sum_{t=1}^{T} \ell_{t}(R^{\star}) \right) + 2\kappa^{2} \log \frac{1}{\delta},$$

which concludes the proof.

D.2.2 PROOF OF LEMMA D.2

Proof of Lemma D.2. Recall the definition of $f: \mathcal{X} \times \mathcal{R} \to \mathbb{R}$ in equation B.4:

$$f(x,R) := -\frac{1}{\eta} \log Z_R(x) + \mathbb{E}_{\pi_R^{\eta}} [R(x,a) - R^{\star}(x,a)].$$

Note f is invariant to adding any action-independent baseline $g: \mathcal{X} \to \mathbb{R}$.

$$f(x, R+g) = -\frac{1}{\eta} \log Z_{R+g}(x) + \mathbb{E}_{\pi_{R+g}^{\eta}} \left[R(x, a) + g(x) - R^{\star}(x, a) \right]$$

$$= -\frac{1}{\eta} \left(\log Z_{R}(x) + \eta g(x) \right) + \mathbb{E}_{\pi_{R}^{\eta}} \left[R(x, a) + g(x) - R^{\star}(x, a) \right] \quad (\pi_{R+g}^{\eta} = \pi_{R}^{\eta})$$

$$= -\frac{1}{\eta} \log Z_{R}(x) + \mathbb{E}_{\pi_{R}^{\eta}} \left[R(x, a) - R^{\star}(x, a) \right] = f(x, R),$$

Algorithm E.1 ODPO (Oracle-efficient Direct Policy Optimization)

- 1: **Inputs:** regularization parameter η , reference policy π_{ref} , online regression oracle OracleLog.
- 2: **Initialize:** choose any $\pi_1 \in \Pi$.
 - 3: **for** round t = 1 to T **do**
 - 4: Observe context $x_t \in \mathcal{X}$.
 - 5: Sample action $a_t^1, a_2^t \sim \pi_t(\cdot|x_t)$ and receive preference feedback y_t .
 - 6: Update π_{t+1} for the next round using OracleDPO via Equation E.2.
 - 7: end for

where the second equality holds because

$$Z_{R+g}(x) = \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|x) e^{\eta(R(x,a) + g(x))} = e^{\eta g(x)} \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|x) e^{\eta R(x,a)} = e^{\eta g(x)} Z_R(x),$$

and

$$\pi_{R+g}^{\eta}(a|x) = \frac{\pi_{\text{ref}}(a|x) \cdot e^{\eta(R(x,a) + g(x))}}{Z_{R+g}(x)} = \frac{\pi_{\text{ref}}(a|x) \cdot e^{\eta R(x,a)} \cdot e^{\eta g(x)}}{e^{\eta g(x)} Z_{R}(x)} = \pi_{R}^{\eta}(a|x).$$

Therefore, by substituting $\hat{R}_t(x,a) \leftarrow \hat{R}_t(x,a) + g(x)$ and the following the proof from Equation B.4 in Lemma B.2, we derive

$$J_t^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_t^{\eta}(\pi_t, R^{\star}) \leqslant \eta \mathbb{E}_{a \sim \pi_t(\cdot | x_t)} \left[\left(\hat{R}_t(x_t, a) - R^{\star}(x_t, a) + g(x_t) \right)^2 \right].$$

which concludes the proof.

E EXTENSION TO DIRECT PREFERENCE OPTIMIZATION (DPO)

In this section, we extend our method to the DPO objective (Rafailov et al., 2023). The problem setup is identical to the RLHF setting (Subsection 3.2), except that DPO bypasses reward learning and directly optimizes the policy within the policy class Π . Rearranging Equation 2, we can express the reward function as follows:

$$R(x,a) = \frac{1}{\eta} \log \frac{\pi(a|x)}{\pi_{\text{ref}}(a|x)} + \frac{1}{\eta} \log Z_R(x).$$
 (E.1)

Accordingly, the Bradley-Terry model for preference feedback takes the form

$$\mathbb{P}(a^1 > a^2 | x, a^1, a^2) = \sigma\left(\frac{1}{\eta} \log \frac{\pi(a^1 | x)}{\pi_{\text{ref}}(a^1 | x)} - \frac{1}{\eta} \log \frac{\pi(a^2 | x)}{\pi_{\text{ref}}(a^2 | x)}\right),\,$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function. Finally, the DPO loss at round t is defined as

$$\ell_t^{\mathrm{DPO}}(\pi) := -\log \sigma \bigg(\frac{1}{\eta} \log \frac{\pi(a_t^1|x_t)}{\pi_{\mathrm{ref}}(a_t^1|x_t)} - \frac{1}{\eta} \log \frac{\pi(a_t^2|x_t)}{\pi_{\mathrm{ref}}(a_t^2|x_t)} \bigg) \,.$$

Note that $\ell_t^{\text{DPO}}(\pi)$ is exactly the same as $\ell_t(R)$ defined in Equation 6.

Similar to Subsection 3.2, we assume access to an online DPO regression oracle, denoted by OracleDPO. At each round t, rather than estimating a reward function, this oracle directly returns a policy:

$$\pi_t \leftarrow \text{OracleDPO}_t((x_t, a_t^1, a_t^2); (x_1, a_1^1, a_1^2, y_1), \dots, (x_{t-1}, a_{t-1}^1, a_{t-1}^2, y_{t-1})).$$
 (E.2)

We assume that the prediction error of OracleDPO is bounded with respect to the policy class Π .

Assumption E.1 (Guarantee of online DPO regression oracle). We assume that, for every (possibly adaptively chosen) sequence $x_{1:T}, a_{1:T}^1, a_{1:T}^2, y_{1:T}$, there exists regret bound $\operatorname{Reg}_{DPO}(T)$ such that the regression oracle OracleDPO satisfies

$$\sum_{t=1}^{T} \ell_t^{\text{DPO}}(\pi_t) - \sum_{t=1}^{T} \ell_t^{\text{DPO}}(\pi_{\eta}^{\star}) \leqslant \text{Reg}_{\text{DPO}}(T).$$

Using this oracle, we establish the following regret bound, analogous to Theorem 3.

Theorem E.1 (Regret of ODPO). Let $\kappa := \sup_{R,x,a} \frac{1}{\dot{\sigma}(R(x,a))}$. Under Assumption E.1, for any $\delta > 0$, ODPO guarantees that with probability at least $1 - \delta$,

$$\begin{split} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) &= \mathcal{O} \left(\eta \kappa^2 \mathrm{Reg}_{\mathrm{DPO}}(T) + \eta \kappa^2 \log(1/\delta) \right), \quad and \\ \mathbf{Regret}(T) &= \mathcal{O} \left(\eta \kappa^2 \mathrm{Reg}_{\mathrm{DPO}}(T) + \eta \kappa^2 \log(1/\delta) + \frac{DT}{\eta} \right), \end{split}$$

where
$$D := \frac{1}{T} \sum_{t=1}^{T} \text{KL}(\pi^{\star}(\cdot || x_t) || \pi_{\text{ref}}(\cdot || x_t)).$$

Proof of Theorem E.1. By Lemma D.1, together with the fact that $\ell_t^{\text{DPO}}(\pi) = \ell_t(R)$ and the reward reformulation in Equation E.1, we obtain

Corollary E.1. With probability at least $1 - \delta$, we have

$$\sum_{t=1}^{T} \left(\frac{1}{\eta} \log \pi_{\eta}^{\star}(a_{t}^{1}|x_{t}) - \frac{1}{\eta} \log \pi_{t}(a_{t}^{1}|x_{t}) - \left(\frac{1}{\eta} \log \pi_{\eta}^{\star}(a_{t}^{2}|x_{t}) - \frac{1}{\eta} \log \pi_{t}(a_{t}^{2}|x_{t}) \right) \right)^{2} \\
\leqslant \kappa^{2} \left(\sum_{t=1}^{T} \ell_{t}^{\text{DPO}}(\pi_{t}) - \sum_{t=1}^{T} \ell_{t}^{\text{DPO}}(\pi_{\eta}^{\star}) \right) + 2\kappa^{2} \log \frac{1}{\delta}.$$

Then, by Lemma D.2, we get

By setting $\delta \leftarrow \frac{\delta}{2}$, we obtain the bound for **Regret**_{KL} (T, η) .

In addition, the bound for $\mathbf{Regret}(T)$ follows directly from Lemma B.3, by applying the same reasoning as in the proof of Theorem 1. This concludes the proof of Theorem E.1.

E.1 Comparison to Lower Bound in Proposition 2.1 of Xie et al. (2024)

A careful reader might wonder whether the logarithmic KL-regularized regret established in Theorem E.1 contradicts the lower bound in Proposition 2.1 of Xie et al. (2024). This is not the case: their analysis considers only the restricted policy class $\Pi = \{\pi_{\rm ref}, \pi_{\eta}^{\star}\}$, rather than the full family of Gibbs policies (Equation 2), so their lower bound does not apply to our setting. For clarity, we first restate Proposition 2.1 from Xie et al. (2024).

Proposition E.1 (Necessity of deliberate exploration, Proposition 2.1 of Xie et al. 2024). Fix $\eta > \frac{8}{\log 2}$, and consider the two-armed bandit setting of $\mathcal{X} = \emptyset$, and $|\mathcal{A}| = N = 2$. Let $\Pi = \{\pi_{\text{ref}}, \pi_{\eta}^{\star}\}$. There exists a reference policy π_{ref} such that for all $T \leq \frac{1}{2} \exp\left(\frac{\eta}{8}\right)$, with constant probability, all of policies π_1, \ldots, π_{T+1} produced by OnlineDPO satisfiy

$$\max_{\pi \in \Pi} J_t^{\eta}(\pi,R) - J_t^{\eta}(\pi_t,R) \geqslant \frac{1}{8}, \quad \forall t \in [T+1].$$

As is clear, this proposition only applies to the restricted class $\Pi = \{\pi_{\rm ref}, \pi_{\eta}^*\}$, where the learner can update its policy only by switching between these two candidates. In contrast, our analysis permits the learner to choose from the full family of Gibbs policies—beyond just $\{\pi_{\rm ref}, \pi_{\eta}^*\}$ —with the choice adaptively guided by data collected through online interactions. Therefore, their lower bound is not directly comparable to our upper bound.

F KL-REGULARIZED CONTEXTUAL BANDITS WITH OFFLINE REGRESSION ORACLE

In this section, we assume access to an *offline regression oracle* instead of the online regression oracle defined in Equation 4. Note that an online regression oracle must provide robust guarantees against arbitrary data sequences generated by an adaptive adversary, which becomes challenging to implement when the function class \mathcal{R} is complex. While the minimax regret rates for online regression with general function classes are well understood (Rakhlin & Sridharan, 2014), to the best of our knowledge, computationally efficient algorithms are only known for specific function classes.

Unlike the online regression oracle setting, where contexts may be chosen adversarially, we now adopt a stochastic context assumption.

Assumption F.1 (Stochastic context). At each round t, the context $x_t \in \mathcal{X}$ is drawn i.i.d. from an unknown but fixed distribution ρ .

In this section, we redefine the *KL-regularized* and *unregularized* regrets in the stochastic contextual setting as follows (we use the same regret notations for simplicity):

$$\begin{split} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) &:= \sum_{t=1}^T \mathbb{E}_{x_t \sim \rho} \big[J_t^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_t^{\eta}(\pi_t, R^{\star}) \big] \quad \text{and} \\ \mathbf{Regret}(T) &:= \sum_{t=1}^T \mathbb{E}_{x_t \sim \rho} \big[\mathbb{E}_{a \sim \pi^{\star}(\cdot \mid x_t)} [R^{\star}(x_t, a)] - \mathbb{E}_{a \sim \pi_t(\cdot \mid x_t)} [R^{\star}(x_t, a)] \big]. \end{split}$$

F.1 OFFLINE REGRESSION ORACLE

We now introduce the notion of an *offline regression oracle*. Given a reward function class \mathcal{R} , an offline regression oracle associated with \mathcal{R} , denoted by OracleOff, is a procedure that produces a predictor $\hat{R}: \mathcal{X} \times \mathcal{A} \to \mathbb{R}$ based on input data. In statistical learning theory, the performance of \hat{R} is typically evaluated in terms of its *out-of-sample error*, that is, its expected error on random, unseen test data. Similar to online regression setting, we assume the statistical learning guarantees of OracleOff.

Assumption F.2 (Guarantee of offline regression oracle). Let $\pi: \mathcal{X} \to \Delta(\mathcal{A})$ be an arbitrary policy. Given n training samples $(x_{1:n}, a_{1:n}, r_{1:n})$ where $x_i \sim \rho$ and $a_i \sim \pi(\cdot \mid x_i)$ i.i.d., the offline regression oracle OracleOff returns a reward estimator $\widehat{R}: \mathcal{X} \times \mathcal{A} \to \mathbb{R}$. For any $\delta > 0$, with probability at least $1 - \delta$, we have

$$\mathbb{E}_{x \sim \rho, a \sim \pi(\cdot | x)} \left[\left(\hat{R}(x, a) - R^{\star}(x, a) \right)^{2} \right] \leqslant \mathcal{E}_{\delta}(n).$$

Under the realizability assumption (Assumption 1), this squared distance corresponds to the estimation error or excess risk of \hat{R} .

Algorithm F.1 KL-EXP-Off

```
1: Inputs: regularization parameter \eta, reference policy \pi_{ref}, offline regression oracle OracleOff,
1730
               epoch schedule 0 = \tau_0 < \tau_1 < \tau_2 < \cdots.
1731
           2: Initialize: choose any \hat{R}_1 \in \mathcal{R}.
1732
           3: for epoch m = 1, 2, ..., m(T) do
1733
                    for round t = \tau_{m-1} + 1, \cdots, \tau_m do
           4:
1734
           5:
                        Observe context x_t \in \mathcal{X}.
1735
                        Compute policy \pi_t(\cdot|x_t) \propto \pi_{\text{ref}}(\cdot|x_t) \exp(\eta R_m(x_t,\cdot)) via Equation 2.
           6:
1736
                         Sample action a_t \sim \pi_t(\cdot|x_t) and receive reward r_t.
           7:
1737
           8:
1738
                    Feed only the data in epoch m-1 into OracleOff and obtain \widehat{R}_{m+1}.
           9:
1739
          10: end for
```

F.2 ALGORITHM AND RESULTS

We provide an algorithm KL-EXP-Off in Algorithm F.1. Unlike Algorithm 1, which updates the predictor at every round, KL-EXP-Off adopts an epoch-based learning protocol, updating the reward estimator only once per epoch via the offline regression oracle. In addition, rather than feeding all past data into the oracle, we restrict its input to the data collected in the immediately preceding epoch (m-1). As a consequence of this strategy, the algorithm proceeds in gradually increasing epochs, i.e., $\tau_m = 2^m$.

Let m(T) denote the total number of epochs. We then establish the following regret bound under the offline regression oracle.

Theorem F.1 (Regret of KL-EXP-Off). Consider an epoch schedule $\tau_m = 2^m$ for $m \leq m(T)$. Then, with probability at least $1 - \delta$, the regret of KL-EXP-Off is bounded by

$$\begin{split} \mathbf{Regret}_{\mathrm{KL}}(T, \eta) &= \mathcal{O} \big(\eta \mathcal{E}_{\delta/\log T}(T) \cdot T \big) \,, \quad \text{and} \\ \mathbf{Regret}(T) &= \mathcal{O} \bigg(\eta \mathcal{E}_{\delta/\log T}(T) \cdot T + \frac{DT}{\eta} \bigg) \,, \end{split}$$

where $D := \frac{1}{T} \sum_{t=1}^{T} \text{KL}(\pi^{\star}(\cdot || x_t) || \pi_{\text{ref}}(\cdot || x_t)).$

Remark F.1 (Computational efficiency). The algorithm KL-EXP-Off requires only $\mathcal{O}(\log T)$ calls to the offline regression oracle.

Example F.1 (Linear classes). When Assumption 1 holds and the reward function class \mathcal{R} is linear (refer Example 1), by using the least squares regression oracle, KL-EXP-Off achieves $\mathbf{Regret}_{\mathrm{KL}}(T,\eta) = \mathcal{O}(\eta d \log T)$ and $\mathbf{Regret}(T) = \mathcal{O}(\sqrt{dDT \log T})$, with the choice $\eta = \Theta\left(\sqrt{\frac{DT}{d \log T}}\right)$. Moreover, by setting π_{ref} to be uniform random, we have $\mathbf{Regret}(T) = \mathcal{O}(\sqrt{dT \log N \log T})$ since $D \leq \log N$. This upper bound matches the lower bound $\Omega(\sqrt{dT \log N \log(T/d)})$ established by Li et al. (2019), up to logarithmic d factors.

Example F.2 (Neural Networks). Let Assumption 1 hold and $\mathcal{R} = \mathcal{G}^N$, where \mathcal{G} denotes the class of Multi-Layer Perceptrons (MLPs) as described in Section 2.1 of Farrell et al. (2021). For each $(x,a) \in \mathcal{X} \times \mathcal{A}$, let the reward function be $R^\star(x,a) = g_a^\star(x)$. Assume the context distribution ρ is continuous over $[-1,1]^d$, and that $g_1^\star,\ldots,g_N^\star$ lie in a Sobolev ball with smoothness $\beta \in \mathbb{N}$. Then, by Theorem 1 of Farrell et al. (2021), the deep MLP-ReLU network estimator attains $\mathcal{O}\left(n^{-\frac{\beta}{\beta+d}}\right)$ estimation error. Consequently, by using this estimator as the offline regression oracle, KL-EXP-Off achieves $\mathbf{Regret}_{\mathrm{KL}}(T,\eta) = \tilde{\mathcal{O}}\left(\eta T^{\frac{d}{\beta+d}}\right)$ and $\mathbf{Regret}(T) = \tilde{\mathcal{O}}\left(T^{\frac{\beta+2d}{2\beta+2d}}\right)$ (ignoring dependence on other parameters) with the parameter choice $\eta = \tilde{\Theta}\left(T^{\frac{\beta}{2\beta+2d}}\right)$. Our derived unregularized regret, $\tilde{\mathcal{O}}\left(T^{\frac{\beta+2d}{2\beta+2d}}\right)$, has the same order as the regret established by Simchi-Levi & Xu (2022).

F.3 Main Proof of Theorem F.1

In this subsection, we present the proof of Theorem F.1.

Proof of Theorem F.1. For any $t \in [T]$, by Lemma B.2, we have

$$\mathbf{Regret}_{\mathrm{KL}}(T,\eta) = \sum_{t=1}^{T} \mathbb{E}_{x_{t} \sim \rho} \left[J_{t}^{\eta}(\pi_{\eta}^{\star}, R^{\star}) - J_{t}^{\eta}(\pi_{t}, R^{\star}) \right]$$

$$\leq \eta \sum_{t=1}^{T} \mathbb{E}_{x_{t} \sim \rho} \mathbb{E}_{a_{t} \sim \pi_{t}(\cdot \mid x_{t})} \left[\left(\hat{R}_{m}(x_{t}, a_{t}) - R^{\star}(x_{t}, a_{t}) \right)^{2} \right] \qquad \text{(Lemma B.2)}$$

Let $\mathcal{F}_t := \sigma(x_1, a_1, r_1, \dots, x_t, r_t, a_t)$ be the filtration up to round t. We introduce the following lemma to further bound the regret.

Lemma F.1 (Lemma 2 of Simchi-Levi & Xu 2022). For all $m \ge 2$ and all $t \in \{\tau_{m-2} + 1, \dots, \tau_{m-1}\}$, with probability at least $1 - \delta/(2m^2)$, we have

$$\mathbb{E}_{x_t \sim \rho, a_t \sim \pi_t(\cdot \mid x_t)} \left[\left(\widehat{R}_m(x_t, a_t) - R^{\star}(x_t, a_t) \right)^2 \mid \mathcal{F}_{t-1} \right] \leqslant \mathcal{E}_{\delta/(2m^2)}(\tau_{m-1} - \tau_{m-2}).$$

By applying Lemma F.1, with probability $1 - \delta$, we obtain

$$\begin{split} \mathbf{Regret}_{\mathrm{KL}}(T,\eta) & \leqslant \eta \sum_{t=1}^{T} \mathbb{E}_{x_{t} \sim \rho} \mathbb{E}_{a_{t} \sim \pi_{t}(\cdot \mid x_{t})} \left[\left(\hat{R}_{m(t)}(x_{t},a_{t}) - R^{\star}(x_{t},a_{t}) \right)^{2} \right] \\ & = \eta \sum_{t=1}^{T} \mathbb{E}_{x_{t} \sim \rho} \mathbb{E}_{a_{t} \sim \pi_{t}(\cdot \mid x_{t})} \left[\left(\hat{R}_{m(t)}(x_{t},a_{t}) - R^{\star}(x_{t},a_{t}) \right)^{2} \mid \mathcal{F}_{t-1} \right] \\ & \leqslant \eta \sum_{t=\tau_{1}+1}^{T} \mathcal{E}_{\delta/(2m(t)^{2})}(\tau_{m(t)-1} - \tau_{m(t)-2}) + \tau_{1} \\ & = \eta \sum_{m=2}^{m(T)} \mathcal{E}_{\delta/(2m^{2})}(\tau_{m-1} - \tau_{m-2}) \cdot (\tau_{m} - \tau_{m-1}) + \tau_{1} \\ & = \mathcal{O}(\eta \mathcal{E}_{\delta/\log T}(T) \cdot T) \; . \end{split}$$

This completes the proof of the upper bound on the KL-regularized regret. Moreover, the bound for the unregularized regret follows directly from the same analysis as in the proof of Theorem 1. \Box

G TECHNICAL LEMMAS

Lemma G.1 (Freedman's inequality, Freedman, 1975). Let $(Z_t)_{t \leq T}$ be a real-valued martingale difference sequence adapted to a filtration \mathcal{F}_{t-1} , and let $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot \mid \mathcal{F}_{t-1}]$. If $|Z_t| \leq B$ almost surely, then for any $\beta \in (0, 1/B)$, it holds that, with probability at least $1 - \delta$,

$$\sum_{t=1}^{T} Z_t \leqslant \beta \sum_{t=1}^{T} \mathbb{E}_{t-1}[Z_t^2] + \frac{B \log(1/\delta)}{\beta}.$$

Lemma G.2 (Lemma A.3 of Foster et al. 2021). Let $(X_t)_{t \leq T}$ be a sequence of random variables adapted to a filtration $(\mathcal{F}_t)_{t \leq T}$. If $0 \leq X_t \leq B$ almost surely, then with probability at least $1 - \delta$,

$$\sum_{t=1}^{T} X_t \leqslant \frac{3}{2} \sum_{t=1}^{T} \mathbb{E}_{t-1}[X_t] + 4B \log \frac{2}{\delta},$$

and

$$\sum_{t=1}^T \mathbb{E}_{t-1}\big[X_t\big] \leqslant 2\sum_{t=1}^T X_t + 8B\log\frac{2}{\delta}.$$

Lemma G.3 (Lemma A.4 of Foster et al. 2021). For any sequence of real-valued random variables $(X_t)_{t \leq T}$ adapted to a filtration $(\mathcal{F}_t)_{t \leq T}$, it holds that with probability at least $1 - \delta$, for all $T' \leq T$,

$$\sum_{t=1}^{T'} X_t \leqslant \sum_{t=1}^{T'} \log \left(\mathbb{E}_{t-1}[e^{X_t}] \right) + \log \frac{1}{\delta}.$$

H ADDITIONAL EXPERIMENTAL RESULTS

H.1 COMPUTATION COST IN LINEAR BANDITS

\overline{N}	$\mid d \mid$	LinUCB	LinTS	LinPHE	SupLinUCB	KL-EXP (ours)
50	5	0.321	0.274	0.862	0.203	0.173
100	5	0.465	0.336	0.927	0.225	0.190
50	20	1.414	1.504	1.877	1.274	1.227
100	20	1.616	1.546	1.942	1.378	1.253

Table H.1: Average per-round computation time (µs) for linear bandits with $d \in \{10, 20\}$ and $N \in \{50, 100\}$.

H.2 COMPUTATION COST IN NEURAL BANDITS

NeuralUCB	NeuralTS	KL-EXP (ours)		
0.0507	0.0665	0.0048		

Table H.2: Average per-round computation time (s) for neural bandits.

H.3 RLHF EXPERIMENTS: DETAILS AND ADDITIONAL RESULTS

In this section, we present the RLHF experimental setup in detail and provide additional results.

Implementation details. For fair comparison, we follow the experimental setup of Dong et al. (2024); Xie et al. (2024). In each iteration, we fix the base model (Llama-3-8B-Flow-SFT) as the reference model $\pi_{\rm ref}$ and set the regularization parameter to $\eta=10.0$. Training is performed with a global batch size of 16, a learning rate of 5×10^{-7} with cosine scheduling, 2 epochs per iteration, and a warmup ratio of 0.03. For XPO, following Xie et al. (2024), we set $\tilde{\pi}_t=\pi_t$ and $\mathcal{D}_t^{\rm opt}=\mathcal{D}_t^{\rm pref}$, and use their exploration schedule $\alpha\in\{1\times 10^{-5},5\times 10^{-6},0\}$ across the three iterations(see their definitions). All experiments were conducted on $8\times$ Nvidia H100 GPUs.

Additional results. Table H.3 reports the accuracies of the algorithms on all 17 academic and chat benchmarks (Zhong et al., 2023; Nie et al., 2019; Hendrycks et al., 2020; Cobbe et al., 2021; Rein et al., 2024; Chen et al., 2021; Zellers et al., 2019; Sakaguchi et al., 2021; Clark et al., 2018; Lin et al., 2021; Mihaylov et al., 2018; Zellers et al., 2018; Sap et al., 2019; Pilehvar & Camacho-Collados, 2018; Levesque et al., 2012; Socher et al., 2013), as well as the performance of OnlineDPO (or ODPO) with varying regularization parameters $\eta \in \{5.0, 8.5, 10.0, 12.5, 20.0\}$. The **bold** values represent the best performance for each benchamrk. The results show that OnlineDPO with a carefully chosen η (= 12.5) outperforms other baselines that rely on additional exploration techniques.

Model	η	iteration	AGIEval	ANLI	MMLU	GSM8K	GPQA	HumanEval	HellaSv	vag Wi	noGrande	ARC-C
Llama-3-8B-Flow-SFT	10.0		39.33	40.51	62.63	74.15	34.34	54.27	59.89)	76.48	53.50
Llama-3-8B-Flow-Final	10.0		41.75	46.29	63.36	74.75	31.31	54.88	61.22	2	76.95	52.73
		iter 1	39.32	43.69	63.08	80.21	34.85	56.10	62.20		75.61	56.40
XP0	10.0	iter 2	39.91	45.77	63.27	80.59	30.81	57.93	62.57		76.32	56.57
		iter 3	40.14	46.07	63.37	82.18	32.83	57.93	62.95	5	76.48	56.57
		iter 1	39.33	45.54	63.10	80.97	33.84	56.10	62.12		76.32	56.14
	5.0	iter 2	40.08	46.13	63.26	82.41	35.35	57.93	62.61		76.24	56.31
	l	iter 3	40.29	46.57	63.15	82.71	33.33	59.75	62.92		76.24	56.66
		iter 1	39.56	45.69	63.22	81.35	31.82	57.32	62.49		76.48	56.31
	8.5	iter 2	40.17	47.19	63.27	82.79	33.33	58.25	62.97		75.85	56.06
	l	iter 3	40.36	48.73	63.31	82.87	32.32	59.73	63.36	5	76.72	55.80
OnlineDPO		iter 1	39.34	44.55	63.25	81.96	32.83	57.93	62.63		76.01	56.14
	10.0	iter 2	40.37	48.13	63.37	82.34	32.83	57.32	63.24		76.16	55.80
	l	iter 3	40.74	48.84	63.25	83.40	32.32	57.93	63.56	5	76.40	56.06
	12.5	iter 1	39.53	45.98	63.26	81.50	31.82	59.76	62.71		76.09	56.06
		iter 2	40.53	47.82	63.30	83.62	31.82	58.54	63.36		76.87	55.12
		iter 3	40.87	48.92	63.30	83.09	33.33	58.54	63.59)	76.32	55.29
	20.0	iter 1	39.72	46.32	63.21	82.87	32.32	56.71	62.98	3	76.01	55.80
		iter 2	40.70	47.80	63.12	82.79	31.82	57.93	63.68		75.93	53.92
		iter 3	41.10	46.29	63.18	83.17	31.31	57.32	63.74	ı	76.40	54.35
Model	η	iteration	ARC-E	Truthful	OA One	nBookOA	SWAG	Social IQa	WiC	WSC273	SST-2	Average
	1 '1	1001011	11110 2									
Liama-3-XB-Flow-SFT	10.0		83 33	45 38	· · ·	35.40						
Llama-3-8B-Flow-Final	10.0		83.33	45.38		35.40 37.20	58.07	52.35	56.74	87.55	90.94	59.11
Llama-3-8B-Flow-Final	10.0	itar 1	81.94	53.71		37.20	58.07 58.15	52.35 52.10	56.74 62.54	87.55 87.18	90.94	59.11
Llama-3-8B-Flow-Final	10.0	iter 1	81.94	53.71 48.53	3	37.20 37.20	58.07 58.15 59.30	52.35 52.10 54.25	56.74 62.54 64.11	87.55 87.18 87.91	90.94 91.97 90.71	59.11 60.47 61.03
		iter 2	81.94 84.09 84.26	53.71 48.53 51.32	3	37.20 37.20 37.80	58.07 58.15 59.30 59.65	52.35 52.10 54.25 53.43	56.74 62.54 64.11 62.54	87.55 87.18 87.91 87.55	90.94 91.97 90.71 90.60	59.11 60.47 61.03 61.23
Llama-3-8B-Flow-Final	10.0	iter 2 iter 3	81.94 84.09 84.26 83.92	53.71 48.53 51.32 52.20	3 2)	37.20 37.20 37.80 38.20	58.07 58.15 59.30 59.65 59.85	52.35 52.10 54.25 53.43 53.12	56.74 62.54 64.11 62.54 61.76	87.55 87.18 87.91 87.55 88.28	90.94 91.97 90.71 90.60 90.60	59.11 60.47 61.03 61.23 61.55
Llama-3-8B-Flow-Final	10.0	iter 2 iter 3	81.94 84.09 84.26 83.92 84.47	53.71 48.53 51.32 52.20 49.41	3 2)	37.20 37.20 37.80 38.20 37.60	58.07 58.15 59.30 59.65 59.85 59.30	52.35 52.10 54.25 53.43 53.12 54.09	56.74 62.54 64.11 62.54 61.76 62.23	87.55 87.18 87.91 87.55 88.28 87.55	90.94 91.97 90.71 90.60 90.60	59.11 60.47 61.03 61.23 61.55 61.09
Llama-3-8B-Flow-Final	10.0	iter 2 iter 3 iter 1 iter 2	81.94 84.09 84.26 83.92 84.47 84.51	53.71 48.53 51.32 52.20 49.41 51.97	3 2)	37.20 37.20 37.80 38.20 37.60 37.20	58.07 58.15 59.30 59.65 59.85 59.30 59.63	52.35 52.10 54.25 53.43 53.12 54.09 53.84	56.74 62.54 64.11 62.54 61.76 62.23 61.44	87.55 87.18 87.91 87.55 88.28 87.55 87.91	90.94 91.97 90.71 90.60 90.60 90.48 90.60	59.11 60.47 61.03 61.23 61.55 61.09 61.61
Llama-3-8B-Flow-Final	10.0	iter 2 iter 3 iter 1 iter 2 iter 3	81.94 84.09 84.26 83.92 84.47 84.51 84.43	53.71 48.53 51.32 52.20 49.41 51.97 53.48	3 2 2))	37.20 37.20 37.80 38.20 37.60 37.20 37.00	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38	90.94 91.97 90.71 90.60 90.48 90.60 91.40	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96
Llama-3-8B-Flow-Final	10.0	iter 2 iter 3 iter 1 iter 2 iter 3 iter 1	81.94 84.09 84.26 83.92 84.47 84.51 84.43 84.22	53.71 48.53 51.32 52.20 49.41 51.97 53.48 51.38	7 3 3 3 3	37.20 37.20 37.80 38.20 37.60 37.20 37.00 37.40	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81 59.53	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96
Llama-3-8B-Flow-Final	10.0	iter 2 iter 3 iter 1 iter 2 iter 3	81.94 84.09 84.26 83.92 84.47 84.51 84.43	53.71 48.53 51.32 52.20 49.41 51.97 53.48	3 2 3 7 3 3	37.20 37.20 37.80 38.20 37.60 37.20 37.00	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38	90.94 91.97 90.71 90.60 90.48 90.60 91.40	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96
Llama-3-8B-Flow-Final XPO	10.0	iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3	81.94 84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09	53.71 48.53 51.32 52.20 49.41 51.97 53.48 51.38 54.42 55.30	3 2 2 3 3 3 3 3 2 2 3 3 3 2 2 3 3 3 2 2 3 3 3 2 2 3 3 3 3 2 3 3 3 2 3	37.20 37.20 37.80 38.20 37.60 37.20 37.00 37.40 37.40 37.60	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81 59.53 59.81 60.00	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.38 52.81	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40 90.71 91.51 90.94	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01
Llama-3-8B-Flow-Final	10.0 10.0 5.0 8.5	iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3	81.94 84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09	53.71 48.53 51.32 52.20 49.41 51.97 53.48 51.38 54.42 55.30	3 2 2 3 3 3 3 2 2 3 3 4 4 4 4 4 4 4 4 4	37.20 37.20 37.80 38.20 37.60 37.20 37.40 37.40 37.40 37.40	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81 59.53 59.81 60.00	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.94 53.38 52.81 53.58	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40 90.71 91.51 90.94 90.83	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01 61.33
Llama-3-8B-Flow-Final XPO	10.0	iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3	81.94 84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09	53.71 48.53 51.32 52.20 49.41 51.97 53.48 51.38 54.42 55.30	3 2 2 3 3 3 2 2 3 3 5 5 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6	37.20 37.20 37.80 38.20 37.60 37.20 37.00 37.40 37.40 37.60	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81 59.53 59.81 60.00	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.38 52.81	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40 90.71 91.51 90.94	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01
Llama-3-8B-Flow-Final XPO	10.0 10.0 5.0 8.5	iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3	84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09 84.55 84.13 83.75	53.71 48.53 51.32 52.20 49.41 51.97 53.48 54.42 55.30 52.01 54.86 56.22	3 2 2 3 3 3 2 3 2 3	37.20 37.20 37.80 38.20 37.60 37.20 37.00 37.40 37.40 37.40 37.40 37.40 37.40	58.07 58.15 59.30 59.65 59.85 59.83 59.63 59.81 60.00 59.50 59.50 59.50	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.38 52.81 53.58 53.12 52.35	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13 61.91 61.60 61.29	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01 88.28 88.64 88.64	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40 90.71 91.51 90.94 90.83 91.06 91.17	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01 61.33 61.77 61.97
Llama-3-8B-Flow-Final XPO	10.0 10.0 5.0 8.5 10.0	iter 2 iter 3 iter 1 iter 2 iter 3	84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09 84.09 84.13 83.75	53.71 48.53 51.32 52.20 49.41 51.97 53.48 51.38 54.42 55.30 52.01 54.86 56.22	3 2 2 3 3 3 3 4 4 4 4 4 4 4 4 4 4 4 4 4	37.20 37.20 37.80 38.20 37.60 37.20 37.00 37.40 37.40 37.40 37.40 37.40 37.40 37.40	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81 60.00 59.50 59.50 59.50 59.51	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.38 52.81 53.58 53.12 52.35 53.07	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13 61.91 61.60 61.29	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01 88.28 88.64 88.64	90.94 91.97 90.71 90.60 90.60 90.60 91.40 90.71 91.51 90.94 90.83 91.06 91.17	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01 61.33 61.77 61.97
Llama-3-8B-Flow-Final XPO	10.0 10.0 5.0 8.5	iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3	84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09 84.55 84.13 83.75	53.71 48.53 51.32 52.20 49.41 51.97 53.48 54.42 55.30 52.01 54.86 56.22	3 2 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3	37.20 37.20 37.80 38.20 37.60 37.20 37.00 37.40 37.40 37.40 37.40 37.40 37.40	58.07 58.15 59.30 59.65 59.85 59.83 59.63 59.81 60.00 59.50 59.50 59.50	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.38 52.81 53.58 53.12 52.35	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13 61.91 61.60 61.29	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01 88.28 88.64 88.64	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40 90.71 91.51 90.94 90.83 91.06 91.17	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01 61.33 61.77 61.97
Llama-3-8B-Flow-Final XPO	10.0 10.0 5.0 8.5 10.0	iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3 iter 1 iter 2 iter 3	84.94 84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09 84.13 83.75 84.39 83.54 83.16	53.71 48.53 51.32 52.20 49.41 51.97 53.48 51.38 54.42 55.30 52.01 54.88 56.22 52.46 55.32 56.70	3 2 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3	37.20 37.80 37.80 38.20 37.60 37.20 37.40 37.40 37.40 37.40 37.40 37.40 37.40 37.40 37.40 37.40 37.40	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81 60.00 59.50 59.91 60.11 59.65 59.95 60.20	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.38 52.81 53.58 53.12 52.35 53.07 52.30 51.84	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13 61.60 61.60 61.60 61.60 61.29	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01 88.28 88.64 88.64 88.64 88.64	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40 90.71 91.51 90.94 90.83 91.06 91.17 90.94 91.63	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01 61.33 61.77 61.97 61.49 61.83 62.09
Llama-3-8B-Flow-Final XPO	10.0 10.0 5.0 8.5 10.0	iter 2 iter 3 iter 1 iter 2 iter 3	84.09 84.26 83.92 84.47 84.51 84.43 84.22 84.09 84.09 84.09 84.55 84.13 83.75 84.39 83.54	53.71 48.53 51.32 52.20 49.41 51.97 53.48 51.38 54.42 55.30 52.01 54.86 56.22 52.46 55.32	3 2 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3	37.20 37.80 37.80 38.20 37.60 37.20 37.00 37.40 37.40 37.40 37.20 37.40 37.40 37.40 37.40 37.40	58.07 58.15 59.30 59.65 59.85 59.30 59.63 59.81 60.00 59.50 59.91 60.11 59.65 59.95	52.35 52.10 54.25 53.43 53.12 54.09 53.84 53.43 53.94 53.38 52.81 53.58 53.12 52.35 53.07 52.30	56.74 62.54 64.11 62.54 61.76 62.23 61.44 62.85 62.54 62.38 61.13 61.91 61.60 61.29 62.07 61.76	87.55 87.18 87.91 87.55 88.28 87.55 87.91 89.38 88.28 87.55 89.01 88.64 88.64 88.64	90.94 91.97 90.71 90.60 90.60 90.48 90.60 91.40 90.71 91.51 90.94 90.83 91.06 91.17	59.11 60.47 61.03 61.23 61.55 61.09 61.61 61.96 61.31 61.79 62.01 61.33 61.77 61.97

Table H.3: Full benchmark evaluation of OnlineDPO with varying $\eta \in \{5.0, 8.5, 10.0, 12.5, 20.0\}$ and of other algorithms that use additional exploration strategies. **Bold** values indicate the best performance.