

# VALUE-INCENTIVIZED PREFERENCE OPTIMIZATION: A UNIFIED APPROACH TO ONLINE AND OFFLINE RLHF

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reinforcement learning from human feedback (RLHF) has demonstrated great promise in aligning large language models (LLMs) with human preference. Depending on the availability of preference data, both online and offline RLHF are active areas of investigation. A key bottleneck is understanding how to incorporate uncertainty estimation in the reward function learned from the preference data for RLHF, regardless of how the preference data is collected. While the principles of optimism or pessimism under uncertainty are well-established in standard reinforcement learning (RL), a practically-implementable and theoretically-grounded form amenable to large language models is not yet available, as standard techniques for constructing confidence intervals become intractable under arbitrary policy parameterizations. In this paper, we introduce a unified approach to online and offline RLHF — value-incentivized preference optimization (VPO) — which regularizes the maximum-likelihood estimate of the reward function with the corresponding value function, modulated by a *sign* to indicate whether the optimism or pessimism is chosen. VPO also directly optimizes the policy with implicit reward modeling, and therefore shares a simpler RLHF pipeline similar to direct preference optimization. Theoretical guarantees of VPO are provided for both online and offline settings, matching the rates of their standard RL counterparts. Moreover, experiments on text summarization, dialogue, and standard benchmarks verify the practicality and effectiveness of VPO.

## 1 INTRODUCTION

Fine-tuning large language models (LLMs) by *reinforcement learning from human feedback* (RLHF) (Ziegler et al., 2019) has been shown to significantly improve the helpfulness, truthfulness and controllability of LLMs, as illustrated by InstructGPT (Ouyang et al., 2022) and many follow-ups. Roughly speaking, there are two critical components of RLHF: (1) *reward modeling*, which maps human preference rankings into a quantitative reward function that can guide policy improvement; and (2) *RL fine-tuning*, which seeks to adjust LLM output to align with human preferences by leveraging the learned reward function, i.e., increasing the probability of preferred answers and decreasing the probability of disfavored answers.

Evidently, the curation of preference data is instrumental in the performance of RLHF, which is commonly modeled as pairwise comparisons from a Bradley-Terry ranking model (Bradley and Terry, 1952). In particular, given a query  $x$ , human annotators choose a preferred answer from two candidate answers  $y_1$  and  $y_2$  generated by an LLM. Despite the simple form, collecting large-scale and high-quality preference data can be expensive and time-consuming. Depending on the availability of preference data, two paradigms of RLHF are considered: (1) *offline* RLHF, where only a pre-collected preference dataset is available, possibly generated from a pre-trained LLM after supervised fine-tuning (SFT); and (2) *online* RLHF, where additional preference data can be collected adaptively to improve alignment. While initial work on RLHF focused on the offline setting, the online setting has also begun to receive considerable attention, as even a small amount of additional preference data has been shown to greatly boost performance.

There has been significant work on the theoretical underpinnings of RLHF that seeks to uncover algorithmic improvements. Notably, while the original RLHF pipeline decouples reward modeling from RL fine-tuning, direct preference optimization (DPO) (Rafailov et al., 2023) integrates these as

a single step in the *offline* setting, leveraging a closed-form solution for the optimal policy in the RL fine-tuning phase. This has led to a welcome simplification of the RLHF pipeline, allowing direct optimization of the policy (i.e., the LLM) from preference data.

Nevertheless, significant challenges remain in RLHF, particularly concerning how to incorporate estimates of reward *uncertainty* in direct preference optimization when parameterizing policies with large-scale neural networks — such as LLMs — in a theoretically and practically effective manner. In standard reinforcement learning (RL), managing uncertainty when an agent interacts with an environment is a critical aspect in achieving near-optimal performance (Sutton and Barto, 2018), when using methods that range from policy-based (Schulman et al., 2017; Xiao et al., 2021), value-based (Mnih et al., 2015; Kumar et al., 2020), and actor-critic methods (Mnih et al., 2016). One dominant approach in the bandit setting, for example, is to construct confidence intervals of the reward estimates, then acting according to the upper and lower confidence bounds — following the principles of optimism and pessimism in the online and offline settings respectively (Lattimore and Szepesvári, 2020; Lai et al., 1985; Rashidinejad et al., 2022).

Despite the fact that uncertainty estimation is even more critical in RLHF, due to the coarse nature of preference data, effective implementations of theoretically justified optimistic and pessimistic principles have yet to be developed in the RLHF literature. For example, existing online preference alignment methods, such as Nash-MD (Munos et al., 2023) and OAIF (Guo et al., 2024), do not incorporate exploration; similarly, pessimism is also not implemented in offline preference alignment methods, such as DPO (Rafailov et al., 2023) and IPO (Azar et al., 2024). A key reason for these omissions is that it is extremely difficult to construct confidence intervals for arbitrary neural networks (Gawlikowski et al., 2021), let alone LLMs. Since optimism for online exploration and pessimism for offline RL both require uncertainty estimation, and given the difficulty of conducting uncertainty estimation for large-scale neural networks, a natural and important question arises:

*Can we implement the optimistic/pessimistic principles under uncertainty in a practically efficient manner for online/offline preference alignment in LLMs while retaining theoretical guarantees?*

## 1.1 OUR CONTRIBUTIONS

In this paper, we provide affirmative answer to the question. Our major contributions are as follows.

- (i) We propose value-incentivized preference optimization (VPO) for both online and offline RLHF, a unified algorithmic framework that *directly optimizes the LLM policy* with the optimistic/pessimistic principles under uncertainty. Avoiding explicit uncertainty estimation, VPO regularizes maximum likelihood estimation of the reward function toward (resp. against) responses that lead to the highest value in the online (resp. offline) setting, hence implementing optimism (resp. pessimism). Theoretical regret guarantees of VPO are developed for both online and offline RLHF, matching their corresponding rates in the standard RL literature with explicit uncertainty estimation.
- (ii) In addition, VPO reveals the critical role of reward calibration, where the shift ambiguity of the reward model inherent in the Bradley-Terry model (Bradley and Terry, 1952) can be exploited to implement additional behavior regularization (Pal et al., 2024; Ethayarajh et al., 2024) via centering the reward model with respect to a *calibration policy*. This allows VPO to provide a theoretical foundation for popular conservative offline RL methods (e.g., (Kumar et al., 2020)), as well as regularized RLHF methods (e.g., DPOP (Pal et al., 2024)).
- (iii) VPO admits a practically-implementable form suitable for RLHF on LLMs, and more generally, deep-learning architectures. We conduct extensive experimental studies using TL;DR and ARC-Challenge tasks as well as standard benchmarks AlpacaEval 2.0 and MT-Bench in online and offline settings with optimistic and pessimistic bias, respectively. The results demonstrate improved empirical performance.

In addition, the proposed value-incentivized regularization can be integrated with other improvements to DPO, e.g., SimPO (Meng et al., 2024) and WPO (Zhou et al., 2024), in a straightforward manner.

## 1.2 RELATED WORK

**RLHF.** Since the introduction of the original RLHF framework, there have been many proposed simplifications of the preference alignment procedure and attempts to improve performance, in-

cluding but not limited to SLiC (Zhao et al., 2023), GSHF (Xiong et al., 2023), DPO (Rafailov et al., 2023), and its variants, such as Nash-MD (Munos et al., 2023), IPO (Azar et al., 2024), OAIF (Guo et al., 2024), SPO (Swamy et al., 2024), SPIN (Chen et al., 2024), GPO (Tang et al., 2024), SimPO (Meng et al., 2024), WPO (Zhou et al., 2024), and DPOP (Pal et al., 2024). These methods can roughly be grouped into online and offline variants, depending on whether preference data is collected before training (offline) or by using the current policy during training (online).

In offline preference alignment, identity preference optimization (IPO, (Azar et al., 2024)) argues that it is problematic to use the Bradley-Terry model in DPO to convert pairwise preferences into pointwise reward values, and proposes an alternative objective function to bypass the use of the Bradley-Terry model. DPO-Positive (DPOP, (Pal et al., 2024)) observes a failure mode of DPO that the standard DPO loss can reduce the model’s likelihood on preferred answers, and proposes to add a regularization term to the DPO objective to avoid such a failure mode. On the other hand, online AI feedback (OAIF, (Guo et al., 2024)) proposes an online version of DPO, where online preference data from LLM annotators is used to evaluate and update the current LLM policy in an iterative manner. Iterative reasoning preference optimization (Iterative RPO, Pang et al. (2024)) proposes to add an additional negative log-likelihood term in the DPO loss to improve performances on reasoning tasks. Finally, Chang et al. (2024) proposes to reuse the offline preference data via reset.

**Uncertainty estimation in RL.** The principles of optimism and pessimism are typically implemented via constructing confidence intervals or posterior sampling, which have been demonstrated to be provably efficient in tabular settings (Jin et al., 2018; Shi et al., 2022). Yet, these approaches have had limited success in conjunction with deep learning architectures (Gawlikowski et al., 2021), and many empirical heuristics in turn lack theoretical validation (Kumar et al., 2020). Notwithstanding, alternative regularization schemes are developed for general function approximation settings with theoretical guarantees, such as Bellman-consistent pessimism for offline RL (Xie et al., 2021), and reward-biased exploration for online RL (Mete et al., 2021; Liu et al., 2024a). VPO draws inspiration from reward-biased exploration (Kumar and Becker, 1982; Liu et al., 2020; 2024a; Hung et al., 2021; Mete et al., 2021) in the standard online RL literature, but significantly broadens its scope to the offline setting and RLHF for the first time.

**Concurrent work.** While preparing this submission, we discovered several concurrent work that also appeared online around the same time proposing similar regularization techniques as ours to encourage optimism (resp. pessimism) for online (resp. offline) RLHF (Zhang et al., 2024a; Xie et al., 2024; Liu et al., 2024b). Despite slightly different forms, the algorithms studied in Zhang et al. (2024a); Xie et al. (2024); Liu et al. (2024b) can be interpreted as adopting different choices of the *calibration policy* in VPO. In the context of online RLHF, Zhang et al. (2024a) empirically studied a similar algorithm as the proposed online VPO under the contextual bandit formulation of RLHF; Xie et al. (2024) provided finite-time regret analysis of a similar algorithm for the token-level MDP formulation with general function approximation, which extends to general deterministic contextual MDPs as well. In the context of offline RLHF, Liu et al. (2024b) studied a similar algorithm as the proposed offline VPO and provided a sample complexity analysis under the contextual bandit formulation, yet focusing on general function approximation and different assumptions.

## 2 PRELIMINARIES

In RLHF, a language model is described by a policy  $\pi$ , which generates an answer  $y \in \mathcal{Y}$  given prompt  $x \in \mathcal{X}$  according to the conditional probability distribution  $\pi(\cdot|x)$ . The standard RLHF process consists of four stages: supervised fine-tuning (SFT), preference data generation, reward modeling, and RL fine-tuning. In the SFT stage, a language model  $\pi_{\text{sft}}$  is obtained by fine-tuning a pre-trained LLM with supervised learning. The remaining stages continue training by leveraging the preference data, which we elaborate below.

**Reward modeling from preference data.** An oracle (e.g., a human labeler or a scoring model) evaluates the quality of two answers  $y_1$  and  $y_2$  given prompt  $x$  and reveals its preference. A widely used approach for modelling the probability of pairwise preferences is the Bradley-Terry model (Bradley and Terry, 1952):

$$\mathbb{P}(y_1 \succ y_2|x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))} = \sigma(r^*(x, y_1) - r^*(x, y_2)), \quad (1)$$

where  $y_1 \succ y_2$  indicates that  $y_1$  is preferred over  $y_2$ ,  $r^* : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is the ground truth reward function, and  $\sigma : \mathbb{R} \rightarrow (0, 1)$  is the logistic function. A preference data sample is denoted by a tuple  $(x, y_+, y_-)$ , where  $y_+$  (resp.  $y_-$ ) is the preferred (resp. unpreferred) answer in the comparison.

Given a preference dataset  $\mathcal{D} = \{(x^i, y_+^i, y_-^i)\}$  composed of independent samples, the reward function  $r$  can be estimated by maximum likelihood estimation (MLE):

$$r_{\text{MLE}} = \arg \min_r \ell(r, \mathcal{D}), \quad (2)$$

where  $\ell(r, \mathcal{D})$  is the negative log-likelihood of  $\mathcal{D}$ , given as  $\ell(r, \mathcal{D}) := - \sum_{(x^i, y_+^i, y_-^i) \in \mathcal{D}} \log \sigma(r(x^i, y_+^i) - r(x^i, y_-^i))$ .

**RL fine-tuning.** Given a reward model  $r$ , we seek to fine-tune the policy  $\pi$  to achieve an ideal balance between the expected reward and its distance from an initial policy  $\pi_{\text{ref}}$ , which is typically the same as  $\pi_{\text{sft}}$ . This is achieved by maximizing the KL-regularized value function  $J(r, \pi)$ , defined as

$$J(r, \pi) = \mathbb{E}_{x \sim \rho, y \sim \pi(\cdot|x)} [r(x, y)] - \beta \mathbb{E}_{x \sim \rho} [\text{KL}(\pi(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x))], \quad (3)$$

where  $\text{KL}(\pi_1 \parallel \pi_2)$  is the KL divergence from  $\pi_1$  to  $\pi_2$ , and  $\beta > 0$  is a regularization parameter. Consequently, the RL fine-tuned policy  $\pi_r$  with respect to the reward  $r$  satisfies  $\pi_r := \arg \max_{\pi} J(r, \pi)$ , which admits a closed-form solution (Rafailov et al., 2023), i.e.,

$$\forall (x \times y) \in \mathcal{X} \times \mathcal{Y} : \quad \pi_r(y|x) = \frac{\pi_{\text{ref}}(y|x) \exp(r(x, y)/\beta)}{Z(r, x)}. \quad (4)$$

Here,  $Z(r, x)$  is a normalization factor given by  $Z(r, x) = \sum_{y' \in \mathcal{Y}} \pi_{\text{ref}}(y'|x) \exp(r(x, y')/\beta)$ .

**Direct preference optimization.** The closed-form solution (4) allows us to write the reward function  $r$  in turn as

$$r(x, y) = \beta(\log \pi_r(y|x) - \log \pi_{\text{ref}}(y|x) + \log Z(r, x)). \quad (5)$$

Plugging the above equation into the reward MLE (2), we obtain the seminal formulation of direct preference optimization (DPO) over the policy space (Rafailov et al., 2023),

$$\pi_{\text{DPO}} = \arg \min_{\pi} - \sum_{(x^i, y_+^i, y_-^i) \in \mathcal{D}} \log \sigma \left( \beta \left( \log \frac{\pi(y_+^i|x)}{\pi_{\text{ref}}(y_+^i|x)} - \log \frac{\pi(y_-^i|x)}{\pi_{\text{ref}}(y_-^i|x)} \right) \right), \quad (6)$$

which avoids explicitly learning the reward model.

### 3 VALUE-INCENTIVIZED PREFERENCE OPTIMIZATION

A major caveat of the standard RLHF framework concerns the lack of accounting for reward uncertainty, which is known to be indispensable in the success of standard RL paradigms in both online and offline settings (Cesa-Bianchi et al., 2017; Rashidinejad et al., 2022). This motivates us to investigate a principled mechanism that be easily integrated into the RLHF pipeline, while bypassing the difficulties of explicit uncertainty estimation in LLMs.

#### 3.1 GENERAL FRAMEWORK

In view of the sub-optimality of naive MLE for reward estimation (Cesa-Bianchi et al., 2017; Rashidinejad et al., 2022), and motivated by the effectiveness of reward-biased MLE in online RL (Kumar and Becker, 1982; Liu et al., 2020; 2024a), we propose to regularize the reward estimate via

$$J^*(r) = \max_{\pi} J(r, \pi), \quad (7)$$

which measures the resulting value function for the given reward if one acts according to its optimal policy. However, in RLHF, by the definition (1), the reward function  $r^*$  is only identifiable up to a prompt-dependent global shift. Specifically, letting  $r_1(x, y) = r_2(x, y) + c(x)$  be two reward functions that only differ by a prompt-dependent shift  $c(x)$ , we have  $r_1(x, y_1) - r_1(x, y_2) = r_2(x, y_1) - r_2(x, y_2)$ , which leads to  $J^*(r_1) = J^*(r_2) + \mathbb{E}_{x \sim \rho}[c(x)]$ . To resolve this challenge, we introduce the following equivalent class of reward functions for the Bradley-Terry model to eliminate the shift ambiguity, which also has the calibration effect of centering the reward function while offering a regularization mechanism to incorporate additional policy preferences.

**Assumption 1** We assume that  $r^* \in \mathcal{R}$ , where

$$\mathcal{R} = \left\{ r : \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [r(x, y)] = 0. \right\}. \quad (8)$$

Here,  $\rho$  is the prompt distribution and  $\pi_{\text{cal}}$  is a fixed calibration distribution independent of the algorithm.

We remark that the number 0 in the condition can be replaced with arbitrary choice of constant, without affecting derivation of our proposed algorithm. This is due to the Bradley-Terry model being invariant to the global shift in the reward function. Therefore, the condition in Assumption 1 only serves the purpose of notational simplicity and does not put a restrict on the reward model in practice.

The proposed regularized MLE of the Bradley-Terry model (2) appends a bias term to the negative likelihood

$$r_{\text{VPO}} = \arg \min_{r \in \mathcal{R}} \{ \ell(r, \mathcal{D}) - \text{sign} \cdot \alpha \cdot J^*(r) \}, \quad (9)$$

incentivizing the algorithm to favor (resp. avoid) reward models with higher value  $J^*(r)$  in the online (resp. offline) setting. Here,  $\alpha > 0$  is a constant controlling the strength of regularization, and  $\text{sign}$  is set to 1 in the online setting and  $-1$  in the offline setting.

At first glance, the objective function for VPO (9) does not immediately imply a computationally-efficient algorithm due to the presence of  $J^*(r)$ . However, by exploiting the same closed-form solution for the optimal policy given the reward in (4), and the reward representation inferred from the policy via (5), we can explicitly express  $J^*(r)$  as

$$\begin{aligned} J^*(r) &= \mathbb{E}_{x \sim \rho, y \sim \pi_r(\cdot|x)} [r(x, y) - \beta(\log \pi_r(y|x) - \log \pi_{\text{ref}}(y|x))] \\ &= \mathbb{E}_{x \sim \rho, y \sim \pi_r(\cdot|x)} [\log Z(r, x)] \\ &= \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\log Z(r, x)] \\ &= \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [r(x, y) - \beta(\log \pi_r(y|x) - \log \pi_{\text{ref}}(y|x))] \\ &= -\beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\log \pi_r(y|x) - \log \pi_{\text{ref}}(y|x)], \end{aligned} \quad (10)$$

where the second step follows because the bracketed term is independent of  $y$  (c.f. (4)) and the last step follows from (8) whenever  $r \in \mathcal{R}$ . Given this key ingredient, we can then rewrite (9) to directly optimize the LLM policy, in a flavor similar to DPO, as

$$\begin{aligned} \pi_{\text{VPO}} &= \operatorname{argmin}_{\pi_r: r \in \mathcal{R}} \{ \ell(r, \mathcal{D}) - \text{sign} \cdot \alpha \cdot J^*(r) \} \\ &= \operatorname{argmin}_{\pi_r: r \in \mathcal{R}} \left\{ - \sum_{(x^i, y_+^i, y_-^i) \in \mathcal{D}} \log \sigma \left( \beta \log \frac{\pi_r(y_+^i|x^i)}{\pi_{\text{ref}}(y_+^i|x^i)} - \beta \log \frac{\pi_r(y_-^i|x^i)}{\pi_{\text{ref}}(y_-^i|x^i)} \right) \right. \\ &\quad \left. + \text{sign} \cdot \alpha \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\log \pi_r(y|x) - \log \pi_{\text{ref}}(y|x)] \right\} \\ &= \operatorname{argmin}_{\pi} \left\{ - \sum_{(x^i, y_+^i, y_-^i) \in \mathcal{D}} \log \sigma \left( \beta \log \frac{\pi(y_+^i|x^i)}{\pi_{\text{ref}}(y_+^i|x^i)} - \beta \log \frac{\pi(y_-^i|x^i)}{\pi_{\text{ref}}(y_-^i|x^i)} \right) \right. \\ &\quad \left. + \text{sign} \cdot \alpha \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\log \pi(y|x) - \log \pi_{\text{ref}}(y|x)] \right\}, \end{aligned} \quad (11)$$

where we drop the constraint on  $r \in \mathcal{R}$ , since for any policy  $\pi$  there exists  $r \in \mathcal{R}$  such that  $\pi = \pi_r$ .

Observing that the reference policy  $\pi_{\text{ref}}(y|x)$  in the last term of (11)  $\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\log \pi(y|x) - \log \pi_{\text{ref}}(y|x)]$  does not impact the optimization solution, we can

change it to  $\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\log \pi(y|x) - \log \pi_{\text{cal}}(y|x)] = - \mathbb{E}_{x \sim \rho} [\text{KL}(\pi_{\text{cal}}(\cdot|x) \parallel \pi(\cdot|x))]$ , which amounts to adding a KL regularization to the original DPO, and offers an interesting interpretation as pushing  $\pi$  against/towards  $\pi_{\text{cal}}$  in the online/offline settings respectively, unveiling the role of reward calibration in RLHF.

**Algorithm 1** VPO for online RLHF

---

**initialization:**  $\pi^{(0)}$ .  
**for**  $t = 0, 1, 2, \dots$  **do**  
    Sample  $x^{(t)} \sim \rho, y_1^{(t)}, y_2^{(t)} \sim \pi^{(t)}(\cdot|x^{(t)})$ .  
    Obtain the preference between  $(x^{(t)}, y_1^{(t)})$  and  $(x^{(t)}, y_2^{(t)})$  from some oracle. Denote the comparison outcome by  $(x^{(t)}, y_+^{(t)}, y_-^{(t)})$ .  
    Update policy  $\pi$  as  

$$\pi^{(t+1)} = \underset{\pi}{\operatorname{argmin}} \left\{ - \sum_{s=1}^t \log \sigma \left( \beta \log \frac{\pi(y_+^{(s)}|x^{(s)})}{\pi_{\text{ref}}(y_+^{(s)}|x^{(s)})} - \beta \log \frac{\pi(y_-^{(s)}|x^{(s)})}{\pi_{\text{ref}}(y_-^{(s)}|x^{(s)})} \right) \right.$$

$$\left. + \alpha \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\log \pi(y|x) - \log \pi_{\text{ref}}(y|x)] \right\}. \quad (14)$$
  
**end**

---

In what follows, we elaborate the development of VPO in both the online and offline settings with corresponding theoretical guarantees under linear function approximation.

### 3.2 ONLINE RLHF: ALGORITHM AND THEORY

The online RLHF procedure extends training by performing reward learning and policy learning iteratively, with a growing preference dataset collected by using the current policy. We use  $\pi^{(t)}$  to denote the policy used in the  $t$ -th iteration, where the superscript  $(t)$  indicates iteration  $t$  in the online setting. The  $t$ -th iteration of VPO for online RLHF consists of the following steps:

- New preference data generation.** We sample a new prompt  $x^{(t)} \sim \rho$  and two answers  $y_1^{(t)}, y_2^{(t)} \sim \pi^{(t)}(\cdot|x^{(t)})$ , query the preference oracle and append  $(x^{(t)}, y_+^{(t)}, y_-^{(t)})$  to the preference dataset.
- Reward learning.** We train a reward model with preference data  $\mathcal{D}^{(t)} := \{(x^{(s)}, y_+^{(s)}, y_-^{(s)})\}_{s=1}^t$  by minimizing the regularized negative log-likelihood, i.e.,

$$r^{(t+1)} = \arg \min_{r \in \mathcal{R}} \{\ell(r, \mathcal{D}^{(t)}) - \alpha \cdot J^*(r)\}. \quad (12)$$

- Policy learning.** This step trains the policy by solving the RL fine-tuning problem:

$$\pi^{(t+1)} = \arg \max_{\pi} J(r^{(t+1)}, \pi). \quad (13)$$

We summarize the detailed procedure in Algorithm 1.

**Theoretical analysis.** Encouragingly, VPO admits appealing theoretical guarantees under function approximation. For simplicity, we restrict attention to linear approximation of the reward model.

**Assumption 2 (Linear Reward)** We parameterize the reward model by

$$r_{\theta}(x, y) = \langle \phi(x, y), \theta \rangle, \quad \forall (x, y) \in \mathcal{X} \times \mathcal{Y}, \quad (15)$$

where  $\phi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$  is a fixed feature mapping and  $\theta \in \mathbb{R}^d$  is the parameters. We assume that  $\|\phi(x, y)\|_2 \leq 1$  for all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ , and that  $r^*(x, y) = \langle \phi(x, y), \theta^* \rangle$  for some  $\theta^*$ .

Under Assumption 1 and 2, it is sufficient to focus on  $\theta \in \Theta$  where

$$\Theta = \left\{ \theta \in \mathbb{R}^d : \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot|x)} [\langle \phi(x, y), \theta \rangle] = 0 \right\}. \quad (16)$$

The next theorem demonstrates that Algorithm 1 achieves  $\tilde{O}(\sqrt{T})$  cumulative regret under mild assumptions. The proof is provided in Appendix B. The proof logic follows from that of (Liu et al., 2024a).

**Theorem 1** Under Assumptions 1 and 2, let  $r_{\theta^{(t)}} \in \Theta$  denote the corresponding reward model for  $\pi^{(t)}$ . Assume that  $\|\theta^*\|_2 \leq C$  and  $\|\theta^{(t)}\|_2 \leq C, \forall t \geq 0$  for some  $C > 0$ . Then with probability  $1 - \delta$  we have

$$\text{Regret} := \sum_{t=1}^T [J^*(r^*) - J(r^*, \pi^{(t)})] \leq \tilde{O}(\exp(2C + C/\beta) \sqrt{\kappa d T}),$$

with  $\alpha = \frac{1}{\exp(2C + C/\beta) \sqrt{\frac{T}{\kappa \min\{d \log T, T\}}}}$  and  $\kappa = \sup_{x, y} \frac{\pi_{\text{cal}}(y|x)}{\pi_{\text{ref}}(y|x)}$ .

**Algorithm 2** VPO for offline RLHF**input:** offline preference data  $\mathcal{D}$  of size  $N$ .Get policy  $\hat{\pi}$  by optimizing

$$\hat{\pi} = \arg \min_{\pi} \left\{ - \sum_{i=1}^N \log \sigma \left( \beta \log \frac{\pi(y_+^i | x^i)}{\pi_{\text{ref}}(y_+^i | x^i)} - \beta \log \frac{\pi(y_-^i | x^i)}{\pi_{\text{ref}}(y_-^i | x^i)} \right) - \alpha \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot | x)} [\log \pi(y | x) - \log \pi_{\text{ref}}(y | x)] \right\}.$$

Theorem 1 shows that VPO achieves the same  $\tilde{O}(\sqrt{T})$  regret for online RLHF as its counterparts in standard contextual bandits with scalar rewards and using UCB for exploration (Lattimore and Szepesvári, 2020).

**Remark 1** The analysis naturally extends to allowing mini-batch samples of size  $M$  in every iteration, yielding an improved regret bound scaled by  $1/\sqrt{M}$  and  $\alpha$  scaled by  $\sqrt{M}$ .

## 3.3 OFFLINE RLHF: ALGORITHM AND THEORY

In offline RLHF, a fixed offline preference dataset is collected  $\mathcal{D} := \{x^i, y_+^i, y_-^i\}_{i=1}^N$ , where  $x^i \sim \rho$ ,  $y^i \sim \pi_b(\cdot | x)$  are sampled from a behavior policy  $\pi_b$ , such as  $\pi_{\text{sft}}$  from SFT. The proposed VPO for offline RLHF consists of one pass through the reward and policy learning phases, i.e.,

$$\hat{r} = \arg \min_{r \in \mathcal{R}} \{ \ell(r, \mathcal{D}) + \alpha \cdot J^*(r) \} \quad \text{and} \quad \hat{\pi} = \arg \max_{\pi} J(\hat{r}, \pi), \quad (17)$$

which discourages over-optimization of the reward function given the limited offline preference data. In the same vein as deriving (14), and by leveraging (10), we obtain the direct policy update rule:

$$\hat{\pi} = \arg \min_{\pi} \left\{ - \sum_{i=1}^N \log \sigma \left( \beta \log \frac{\pi(y_+^i | x^i)}{\pi_{\text{ref}}(y_+^i | x^i)} - \beta \log \frac{\pi(y_-^i | x^i)}{\pi_{\text{ref}}(y_-^i | x^i)} \right) - \alpha \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}(\cdot | x)} [\log \pi(y | x) - \log \pi_{\text{ref}}(y | x)] \right\}. \quad (18)$$

We summarize the detailed procedure in Algorithm 2. When  $\pi_{\text{cal}}$  is set to  $\pi_{\text{ref}}$ , the regularization term becomes the KL divergence between  $\pi$  and  $\pi_{\text{ref}}$ , which is reminiscent of a popular choice in offline RL practice (Kumar et al., 2020). Another heuristic choice is to set  $\pi_{\text{cal}}$  to the marginalized positive answer distribution from the dataset, i.e.,  $(x, y_+) \sim \mathcal{D}$ , which leads to a similar objective in (Pal et al., 2024).

**Saddle-point characterization and pessimism.** We first illustrate that VPO indeed executes the principle of pessimism in a complementary manner to the standard approach of pessimism, which finds a policy that maximizes the worst-case value function over a confidence set. In particular, this strategy (Uehara and Sun, 2021) obtains a policy by solving

$$\hat{\pi}_{\text{LCB}} = \arg \max_{\pi} \min_{r \in \mathcal{R}_{\delta}} J(r, \pi) \quad (19)$$

where the confidence set  $\mathcal{R}_{\delta}$  is typically set to  $\{r : \ell(r, \mathcal{D}) \leq \ell(r_{\text{MLE}}, \mathcal{D}) + \delta\}$  or  $\{r : \text{dist}(r, r_{\text{MLE}}) \leq \delta\}$  for some  $\delta > 0$  and s distance measure  $\text{dist}$ . Turning to VPO, note that by (17) we have

$$\hat{r} = \arg \min_r \{ \ell(r, \mathcal{D}) + \alpha J^*(r) \} = \arg \min_r \max_{\pi} \{ \ell(r, \mathcal{D}) + \alpha J(r, \pi) \}. \quad (20)$$

Since  $\ell(r, \mathcal{D}) + \alpha J(r, \pi)$  is strongly concave over  $\pi$ , and convex over  $r$ , it allows us to formulate  $(\hat{r}, \hat{\pi})$  as a saddle point in the following lemma. The proof is given in Appendix C.1.

**Lemma 1**  $(\hat{r}, \hat{\pi})$  is a saddle point of the objective  $\ell(r, \mathcal{D}) + \alpha J(r, \pi)$ , i.e., for any  $(r', \pi')$ , we have

$$\begin{cases} \ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi}) \leq \ell(r', \mathcal{D}) + \alpha J(r', \hat{\pi}) \\ \ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi}) \geq \ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \pi') \end{cases}.$$

As such, the policy obtained by VPO can be equivalently written as

$$\hat{\pi} \in \arg \max_{\pi} \min_r \left\{ J(r, \pi) + \frac{1}{\alpha} \ell(r, \mathcal{D}) \right\} = \arg \max_{\pi} \min_{r \in \mathcal{R}_{\delta}(\pi, \alpha)} J(r, \pi), \quad (21)$$

where  $\mathcal{R}_{\delta}(\pi, \alpha)$  is the constraint set  $\{r : \ell(r, \mathcal{D}) \leq \ell(r_{\text{MLE}}, \mathcal{D}) + \delta(\pi, \alpha)\}$  such that the constrained optimization problem  $\min_{r \in \mathcal{R}_{\delta}(\pi, \alpha)} J(r, \pi)$  is equivalent to the regularized problem  $\min_r \{ J(r, \pi) + \frac{1}{\alpha} \ell(r, \mathcal{D}) \}$ . In view of the similarity between the formulations (19) and (21), we conclude that VPO

implements the pessimism principle (19) in an oblivious manner without explicitly estimating the uncertainty level, justifying popular practice as a valid approach to pessimism (Kumar et al., 2020).

**Theoretical analysis.** The next theorem establishes the sub-optimality gap of VPO with linear function approximation under mild assumptions. The proof is given in Appendix C.

**Theorem 2** *Under Assumptions 1 and 2, let  $\hat{\theta} \in \Theta$  denote the corresponding reward model for  $\hat{\pi}$ . Assume that  $\|\theta^*\|_2 \leq C$  and  $\|\hat{\theta}\|_2 \leq C$  for some  $C > 0$ . Let  $\alpha = \sqrt{N}$  and  $\delta \in (0, 1)$ . With probability  $1 - \delta$ , we have*

$$J^*(r^*) - J(r^*, \hat{\pi}) \leq \mathcal{O}\left(\frac{C_1}{\sqrt{N}} \cdot \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} + \frac{C_2}{\sqrt{N}}\right),$$

where  $\Sigma_{\mathcal{D}} = \frac{1}{N} \sum_{i=1}^N (\phi(x^i, y_+^i) - \phi(x^i, y_-^i))(\phi(x^i, y_+^i) - \phi(x^i, y_-^i))^\top$  is the feature sample covariance matrix,  $\lambda = 1/N$ ,  $C_1 = \exp(C) \left( \sqrt{d + \log(1/\delta)} + \kappa_{\mathcal{D}} \right) + C$  and  $C_2 = \exp(C) \kappa_{\mathcal{D}}^2 + C \kappa_{\mathcal{D}} + 1$ . Here,

$$\kappa_{\mathcal{D}} = \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} \leq 4(\lambda_{\min}(\Sigma_{\mathcal{D}}) + \lambda)^{-1}.$$

Theorem 2 establishes that VPO achieves the same rate of  $\tilde{\mathcal{O}}(1/\sqrt{N})$  as standard offline RL, as long as the offline dataset  $\mathcal{D}$  has sufficient coverage. We remark that  $\left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}}$  is

reminiscent of the standard single-policy concentratability coefficient in offline RL, which measures the distribution shift between the offline dataset and the optimal policy (Zhu et al., 2023).

**Remark 2** *Recently, Rafailov et al. (2024); Zhong et al. (2024) offered an interpretation of DPO using the token-level Markov Decision Process (MDP), aiming at reconciling the gap between the practical fine-tuning of LLMs at the token level and the theoretical formulation of DPO at the sentence level. Fortunately, VPO can be interpreted using the same setup without introducing further algorithmic modifications. We provide the detailed discussion in Appendix A.*

## 4 EXPERIMENTS

In this section, We evaluate the pessimistic/optimistic VPO for LLMs in offline and online setting, respectively. In both settings, the proposed VPO demonstrates strong performances over the baselines. We provide additional experiments on synthetic bandits in Appendix D.3.

### 4.1 OFFLINE SETTING

In this setting, we focus on ARC-Challenge task (Clark et al., 2018), which consists of multiple-choice questions across various science subjects, each with a ground truth answer. To construct the preference pairs for training, we start with 1, 119 examples in the training set and generate three comparison pairs with each incorrect answer, resulting in a total of 3, 357 preference training data. We use ARC-Challenge test set which contains 1, 172 questions to test algorithms performances.

We evaluate pessimistic VPO and compare its performance to several offline RLHF baselines (DPO (Rafailov et al., 2023) and IPO (Azar et al., 2024)) on several LLMs, including LLAMA2-7B-CHAT, LLAMA2-13B-CHAT (Touvron et al., 2023) and FLAN-T5-XL (Chung et al., 2022). We emphasize that our goal is to evaluate the RLHF algorithm designs for LLMs, rather than pushing LLM towards state-of-the-art performance. For fair comparison, we keep all the experiment settings and prompts the same for every RLHF algorithm. We did not apply any additional chain-of-thought reasoning to avoid compounding factors affecting the RLHF performances. We tuned the hyperparameters for VPO and the baselines on the validation set to achieve their best performances. For detailed hyperparameters setup and prompting strategy, please refer to Appendix D.

The performances are illustrated in Figure 1. As we can see, the proposed VPO method demonstrates significantly better performance compared to IPO. Another important observation is that the proposed VPO method is more robust to over-optimization (Gao et al., 2023) compared to DPO. As DPO training continues, its performance declines. In contrast, VPO consistently maintains the performances, avoiding the over-optimization issue and justifying the implicit robustness of pessimism as we revealed in (20).



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

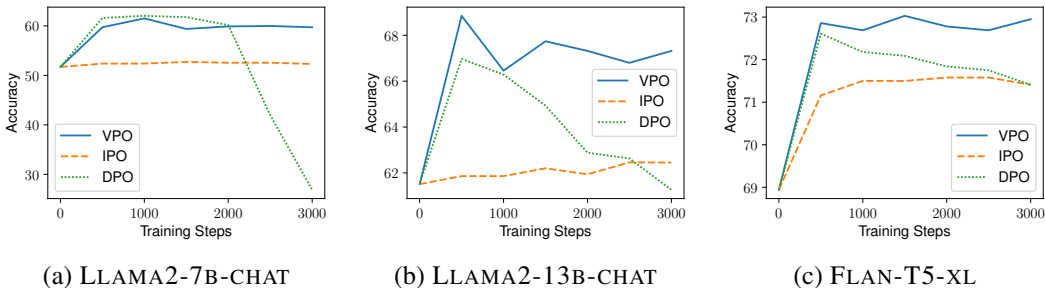


Figure 1: The accuracy of the LLAMA2-7B-CHAT, LLAMA2-13B-CHAT and FLAN-T5-XL policies trained by VPO and other baselines (DPO and IPO) on ARC-challenge, respectively. The proposed pessimistic VPO performs consistently strong, and avoids over-optimization.

4.2 ONLINE SETTING

For online setting, we conduct two distinct experimental setups. The first, referred to as **Buffer**, adopts the experimental setup in Online AI Feedback (OAIF) (Guo et al., 2024), with an additional buffer used to sample the data for the exploration part of the VPO loss. The second setup, referred to as **Iterative**, adopts the experimental setup in (Zhang et al., 2024a), relying on an online iterative training framework.

**Buffer.** In these set of experiments, we adopt OAIF experimental setup (Guo et al., 2024) where the preference data is gathered by online sampling from the policy and labeled through online feedback. We also introduce a buffer that stores the labeled preferences and is used to sample the data for the exploration term in the VPO loss. We adopt PALM2-XXS language model (Anil et al., 2023) as policy, initialized by supervised finetuning, denoted as SFT model. We exploit another PALM2-XS model as the LLM annotator to provide online feedback. We evaluate the performance of optimistic VPO and compare its performance to Online DPO (Guo et al., 2024). We choose TL;DR task (Stiennon et al., 2020) and extract its prompts for the input of preference data. Similar to (Guo et al., 2024), we use *Detailed 0-shot* prompt from Lee et al. (2023). The prompts we used and how we get preference scores are detailed in Appendix D. We emphasize our algorithm is agnostic to human or AI feedback.

As a sanity check, we track the win rate of VPO and Online DPO against the SFT baseline on TL;DR during training in Figure 2a. For ablation purpose, we vary the exploration weight  $\alpha = \{0.01, 0.1\}$  in the optimistic VPO. One significant observation is that although all the online RLHF algorithms follow the increase trend, the win-rate against SFT of the optimistic VPO has larger oscillation, comparing to Online DPO. And the oscillation reduces, with  $\alpha$  diminishing. Our conjecture is that this behavior is encouraged by the optimistic term in VPO, for collecting more unexplored data, which may delay the learning due to the diversity in data. However, as the learning proceeds, the proposed VPO outperforms the competitors, because of the coverage of the collected data. To

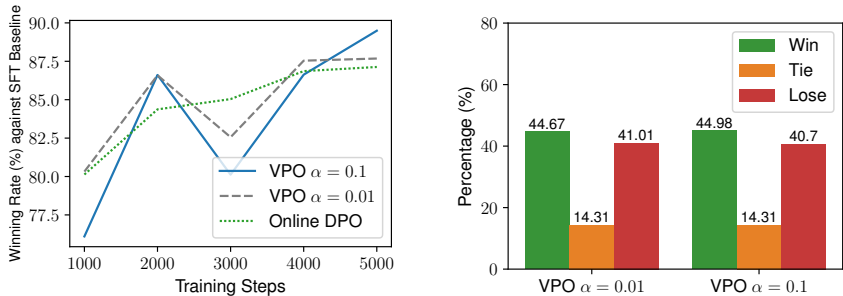


Figure 2(a): Win rate of VPO and Online DPO against the SFT baseline on TL;DR task.

Figure 2(b): Win/tie/loss rate of VPO with different exploration rate  $\alpha = \{0.01, 0.1\}$ , directly against Online DPO.

demonstrate the advantages of optimistic VPO in online setting more directly, we evaluate the win/tie/loss rate against Online DPO head-to-head, as shown in Figure 2b. This clearly shows

that the optimistic VPO achieves better performances with larger exploration preference, and thus, consolidates our conclusion that *i*), the simple value-incentivized term makes the exploration practical without uncertainty estimation; and *ii*), exploration is potentially beneficial for better model.

**Iterative.** We further evaluate the performance of VPO on standard benchmarks AlpacaEval 2.0 (Dubois et al., 2024b;a) and MT-Bench (Zheng et al., 2024), using similar experimental setup in recent literature (Zhang et al., 2024b; Wu et al., 2024). We use UltraFeedback<sup>1</sup> (Cui et al., 2023) as our training dataset which contains around 61k preference pairs of single-turn conversations. We split the 61k prompts into four chunks and follow an iterative training approach. We choose Zephyr-7B-SFT<sup>2</sup> as our LLM model. We follow the best hyperparameters setup found in Zhang et al. (2024b) and compare the results of VPO to DPO reported therein. We first conduct a single iteration of standard DPO training on the first portion of the training data, referred to as Zephyr-7B-DPO in 1. We then perform 3 iterations of VPO, each iteration on a different data portion, while using online AI feedback provided by PairRM (Jiang et al., 2023) in between. Further details of our experiments are explained in Appendix D.

The evaluation results are summarized in Table 1, filling out the baselines based on Zhang et al. (2024b). As could be seen, VPO significantly improves the performance of the base model Zephyr-7B-SFT by **14.52**, achieving the highest length-controlled (LC) Win Rate on the AlpacaEval 2.0 benchmark, beating DPO. This result, 22.53, is competitive to much larger models such as much Yi-34B-Chat, 27.19, and Llama-3-70B-Instruct, 33.17. Additionally, VPO shows significant improvement on MT-Bench compared to the base model Zephyr-7B-SFT with the increase of **2.32** while beating DPO. We further report the results of other iterative post-training algorithms, such as SPIN (Chen et al., 2024), DNO (Rosset et al., 2024), and SPPO (Wu et al., 2024) and show that even though VPO is trained on a weak base model, it achieves comparable results to these baselines. Granular views on a radar chart can be found in Appendix D.

Model	AlpacaEval 2.0			MT-Bench		
	LC Win Rate	Win Rate	Avg. Len	Avg	1st Turn	2nd Turn
Zephyr-7B-SFT	8.01	4.63	916	5.30	5.63	4.97
Zephyr-7B-DPO	15.41	14.44	1752	7.31	7.55	7.07
DPO Iter 1 (Zephyr)	20.53	16.69	1598	7.53	7.81	7.25
DPO Iter 2 (Zephyr)	22.12	19.82	1717	7.55	7.85	7.24
DPO Iter 3 (Zephyr)	22.19	<b>19.88</b>	1717	7.46	7.85	7.06
VPO Iter 1 (Zephyr)	<b>22.53</b>	19.09	1638	7.50	7.76	7.24
VPO Iter 2 (Zephyr)	21.84	18.78	1663	<b>7.62</b>	7.93	<b>7.32</b>
VPO Iter 3 (Zephyr)	22.15	19.58	1713	7.61	<b>8.01</b>	7.21
SPIN	7.23	6.54	1426	6.54	6.94	6.14
Orca-2.5-SFT	10.76	6.99	1174	6.88	7.72	6.02
DNO (Orca-2.5-SFT)	22.59	24.97	2228	7.48	7.62	7.35
Mistral-7B-Instruct-v0.2	19.39	15.75	1565	7.51	7.78	7.25
SPPO (Mistral-it)	28.53	31.02	2163	7.59	7.84	7.34
Yi-34B-Chat	27.19	21.23	2123	7.90	-	-
Llama-3-70B-Instruct	33.17	33.18	1919	9.01	9.21	8.80
GPT-4 Turbo (04/09)	55.02	46.12	1802	9.19	9.38	9.00

Table 1: Results on AlpacaEval 2.0 and MT-Bench.

## 5 CONCLUSION AND DISCUSSION

In this work, we develop a unified approach to achieving principled optimism and pessimism in online and offline RLHF, which enables a practical computation scheme by incorporating uncertainty estimation implicitly within reward-biased maximum likelihood estimation. Theoretical analysis indicates that the proposed methods mirror the guarantees of their standard RL counterparts, which is furthermore corroborated by numerical results. Important future directions include investigating adaptive rules for selecting  $\alpha$  without prior information and more refined analysis on the choice of  $\pi_{\text{cal}}$ . This work also hints at a general methodology of designing practical algorithms with principled optimism/pessimism under more general RL setups.

<sup>1</sup>[https://huggingface.co/datasets/HuggingFaceH4/ultrafeedback\\_binarized](https://huggingface.co/datasets/HuggingFaceH4/ultrafeedback_binarized)

<sup>2</sup><https://huggingface.co/HuggingFaceH4/mistral-7b-sft-beta>

## REFERENCES

- 540  
541  
542 Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. (2011). Improved algorithms for linear stochastic  
543 bandits. *Advances in neural information processing systems*, 24.
- 544 Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey,  
545 P., Chen, Z., et al. (2023). Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- 546  
547 Azar, M. G., Guo, Z. D., Piot, B., Munos, R., Rowland, M., Valko, M., and Calandriello, D. (2024).  
548 A general theoretical paradigm to understand learning from human preferences. In *International*  
549 *Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR.
- 550 Bradley, R. A. and Terry, M. E. (1952). Rank analysis of incomplete block designs: I. the method of  
551 paired comparisons. *Biometrika*, 39(3/4):324–345.
- 552  
553 Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. (2022). Fast global convergence of natural policy  
554 gradient methods with entropy regularization. *Operations Research*, 70(4):2563–2578.
- 555 Cesa-Bianchi, N., Gentile, C., Lugosi, G., and Neu, G. (2017). Boltzmann exploration done right.  
556 *Advances in neural information processing systems*, 30.
- 557  
558 Chang, J. D., Shan, W., Oertell, O., Brantley, K., Misra, D., Lee, J. D., and Sun, W. (2024). Dataset  
559 reset policy optimization for RLHF. *arXiv preprint arXiv:2404.08495*.
- 560 Chen, Z., Deng, Y., Yuan, H., Ji, K., and Gu, Q. (2024). Self-play fine-tuning converts weak language  
561 models to strong language models. *arXiv preprint arXiv:2401.01335*.
- 562  
563 Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, E., Wang, X., Dehghani, M.,  
564 Brahma, S., et al. (2022). H. chi, jeff dean, jacob devlin, adam roberts, denny zhou, quoc v. le, and  
565 jason wei. 2022. scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- 566 Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. (2018).  
567 Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint*  
568 *arXiv:1803.05457*.
- 569 Cui, G., Yuan, L., Ding, N., Yao, G., Zhu, W., Ni, Y., Xie, G., Liu, Z., and Sun, M. (2023). Ultrafeed-  
570 back: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*.
- 571  
572 Dubois, Y., Galambosi, B., Liang, P., and Hashimoto, T. B. (2024a). Length-controlled alpacaeval: A  
573 simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.
- 574  
575 Dubois, Y., Li, C. X., Taori, R., Zhang, T., Gulrajani, I., Ba, J., Guestrin, C., Liang, P. S., and  
576 Hashimoto, T. B. (2024b). AlpacaFarm: A simulation framework for methods that learn from  
577 human feedback. *Advances in Neural Information Processing Systems*, 36.
- 578  
579 Ethayarajh, K., Xu, W., Muennighoff, N., Jurafsky, D., and Kiela, D. (2024). KTO: Model alignment  
as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*.
- 580  
581 Gao, L., Schulman, J., and Hilton, J. (2023). Scaling laws for reward model overoptimization. In  
582 *International Conference on Machine Learning*, pages 10835–10866. PMLR.
- 583  
584 Gawlikowski, J., Tassi, C. R. N., Ali, M., Lee, J., Humt, M., Feng, J., Kruspe, A., Triebel, R., Jung,  
585 P., Roscher, R., et al. (2021). A survey of uncertainty in deep neural networks. *arXiv preprint*  
*arXiv:2107.03342*.
- 586  
587 Guo, S., Zhang, B., Liu, T., Liu, T., Khalman, M., Llinares, F., Rame, A., Mesnard, T., Zhao, Y.,  
588 Piot, B., et al. (2024). Direct language model alignment from online ai feedback. *arXiv preprint*  
*arXiv:2402.04792*.
- 589  
590 Hung, Y.-H., Hsieh, P.-C., Liu, X., and Kumar, P. (2021). Reward-biased maximum likelihood  
591 estimation for linear stochastic bandits. In *Proceedings of the AAAI Conference on Artificial*  
592 *Intelligence*, volume 35, pages 7874–7882.
- 593  
Jiang, D., Ren, X., and Lin, B. Y. (2023). Llm-blender: Ensembling large language models with  
pairwise ranking and generative fusion. *arXiv preprint arXiv:2306.02561*.

- 594 Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is Q-learning provably efficient?  
595 *Advances in neural information processing systems*, 31.
- 596
- 597 Jin, C., Liu, Q., and Yu, T. (2022). The power of exploiter: Provable multi-agent rl in large state  
598 spaces. In *International Conference on Machine Learning*, pages 10251–10279. PMLR.
- 599
- 600 Kumar, A., Zhou, A., Tucker, G., and Levine, S. (2020). Conservative q-learning for offline  
601 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191.
- 602
- 603 Kumar, P. and Becker, A. (1982). A new family of optimal adaptive controllers for markov chains.  
604 *IEEE Transactions on Automatic Control*, 27(1):137–146.
- 605
- 606 Lai, T. L., Robbins, H., et al. (1985). Asymptotically efficient adaptive allocation rules. *Advances in  
607 applied mathematics*, 6(1):4–22.
- 608
- 609 Lattimore, T. and Szepesvári, C. (2020). *Bandit algorithms*. Cambridge University Press.
- 610
- 611 Lee, H., Phatale, S., Mansoor, H., Lu, K., Mesnard, T., Bishop, C., Carbune, V., and Rastogi, A.  
612 (2023). Rlaif: Scaling reinforcement learning from human feedback with ai feedback. *arXiv  
613 preprint arXiv:2309.00267*.
- 614
- 615 Liu, X., Hsieh, P.-C., Hung, Y. H., Bhattacharya, A., and Kumar, P. (2020). Exploration through  
616 reward biasing: Reward-biased maximum likelihood estimation for stochastic multi-armed bandits.  
617 In *International Conference on Machine Learning*, pages 6248–6258. PMLR.
- 618
- 619 Liu, Z., Lu, M., Xiong, W., Zhong, H., Hu, H., Zhang, S., Zheng, S., Yang, Z., and Wang, Z.  
620 (2024a). Maximize to explore: One objective function fusing estimation, planning, and exploration.  
621 *Advances in Neural Information Processing Systems*, 36.
- 622
- 623 Liu, Z., Lu, M., Zhang, S., Liu, B., Guo, H., Yang, Y., Blanchet, J., and Wang, Z. (2024b). Provably  
624 mitigating overoptimization in RLHF: Your SFT loss is implicitly an adversarial regularizer. *arXiv  
625 preprint arXiv:2405.16436*.
- 626
- 627 Meng, Y., Xia, M., and Chen, D. (2024). Simpo: Simple preference optimization with a reference-free  
628 reward. *arXiv preprint arXiv:2405.14734*.
- 629
- 630 Mete, A., Singh, R., Liu, X., and Kumar, P. (2021). Reward biased maximum likelihood estimation  
631 for reinforcement learning. In *Learning for Dynamics and Control*, pages 815–827. PMLR.
- 632
- 633 Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu,  
634 K. (2016). Asynchronous methods for deep reinforcement learning. In *International conference  
635 on machine learning*, pages 1928–1937. PMLR.
- 636
- 637 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A.,  
638 Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep  
639 reinforcement learning. *nature*, 518(7540):529–533.
- 640
- 641 Munos, R., Valko, M., Calandriello, D., Azar, M. G., Rowland, M., Guo, Z. D., Tang, Y., Geist,  
642 M., Mesnard, T., Michi, A., et al. (2023). Nash learning from human feedback. *arXiv preprint  
643 arXiv:2312.00886*.
- 644
- 645 Nachum, O., Norouzi, M., Xu, K., and Schuurmans, D. (2017). Bridging the gap between value and  
646 policy based reinforcement learning. *Advances in neural information processing systems*, 30.
- 647
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S.,  
Slama, K., Ray, A., et al. (2022). Training language models to follow instructions with human  
feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Pal, A., Karkhanis, D., Dooley, S., Roberts, M., Naidu, S., and White, C. (2024). Smaug: Fixing  
failure modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*.
- Pang, R. Y., Yuan, W., Cho, K., He, H., Sukhbaatar, S., and Weston, J. (2024). Iterative reasoning  
preference optimization. *arXiv preprint arXiv:2404.19733*.

- 648 Rafailov, R., Hejna, J., Park, R., and Finn, C. (2024). From  $r$  to  $q^*$ : Your language model is secretly  
649 a Q-function. *arXiv preprint arXiv:2404.12358*.  
650
- 651 Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. (2023). Direct  
652 preference optimization: Your language model is secretly a reward model. *Advances in Neural  
653 Information Processing Systems*, 36.
- 654 Rashidinejad, P., Zhu, B., Ma, C., Jiao, J., and Russell, S. (2022). Bridging offline reinforcement  
655 learning and imitation learning: A tale of pessimism. *IEEE Transactions on Information Theory*,  
656 68(12):8156–8196.  
657
- 658 Rosset, C., Cheng, C.-A., Mitra, A., Santacrose, M., Awadallah, A., and Xie, T. (2024). Direct nash  
659 optimization: Teaching language models to self-improve with general preferences. *arXiv preprint  
660 arXiv:2404.03715*.
- 661 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy  
662 optimization algorithms. *arXiv preprint arXiv:1707.06347*.  
663
- 664 Shi, L., Li, G., Wei, Y., Chen, Y., and Chi, Y. (2022). Pessimistic Q-learning for offline reinforcement  
665 learning: Towards optimal sample complexity. In *International conference on machine learning*,  
666 pages 19967–20025. PMLR.
- 667 Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and  
668 Christiano, P. F. (2020). Learning to summarize with human feedback. *Advances in Neural  
669 Information Processing Systems*, 33:3008–3021.  
670
- 671 Sutton, R. S. and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press.  
672
- 673 Swamy, G., Dann, C., Kidambi, R., Wu, Z. S., and Agarwal, A. (2024). A minimaximalist approach  
674 to reinforcement learning from human feedback. *arXiv preprint arXiv:2401.04056*.
- 675 Tang, Y., Guo, Z. D., Zheng, Z., Calandriello, D., Munos, R., Rowland, M., Richemond, P. H., Valko,  
676 M., Pires, B. Á., and Piot, B. (2024). Generalized preference optimization: A unified approach to  
677 offline alignment. *arXiv preprint arXiv:2402.05749*.  
678
- 679 Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S.,  
680 Bhargava, P., Bhosale, S., et al. (2023). Llama 2: Open foundation and fine-tuned chat models.  
681 *arXiv preprint arXiv:2307.09288*.
- 682 Uehara, M. and Sun, W. (2021). Pessimistic model-based offline reinforcement learning under partial  
683 coverage. *arXiv preprint arXiv:2107.06226*.  
684
- 685 Wu, Y., Sun, Z., Yuan, H., Ji, K., Yang, Y., and Gu, Q. (2024). Self-play preference optimization for  
686 language model alignment. *arXiv preprint arXiv:2405.00675*.  
687
- 688 Xiao, C., Wu, Y., Mei, J., Dai, B., Lattimore, T., Li, L., Szepesvari, C., and Schuurmans, D. (2021).  
689 On the optimality of batch policy optimization algorithms. In *International Conference on Machine  
690 Learning*, pages 11362–11371. PMLR.
- 691 Xie, T., Cheng, C.-A., Jiang, N., Mineiro, P., and Agarwal, A. (2021). Bellman-consistent pessimism  
692 for offline reinforcement learning. *Advances in neural information processing systems*, 34:6683–  
693 6694.
- 694 Xie, T., Foster, D. J., Krishnamurthy, A., Rosset, C., Awadallah, A., and Rakhlin, A. (2024).  
695 Exploratory preference optimization: Harnessing implicit  $q^*$ -approximation for sample-efficient  
696 rlhf.  
697
- 698 Xiong, W., Dong, H., Ye, C., Zhong, H., Jiang, N., and Zhang, T. (2023). Gibbs sampling from human  
699 feedback: A provable KL-constrained framework for RLHF. *arXiv preprint arXiv:2312.11456*.  
700
- 701 Zhan, W., Uehara, M., Kallus, N., Lee, J. D., and Sun, W. (2023). Provable offline preference-based  
reinforcement learning. In *The Twelfth International Conference on Learning Representations*.

- Zhang, S., Yu, D., Sharma, H., Yang, Z., Wang, S., Hassan, H., and Wang, Z. (2024a). Self-exploring language models: Active preference elicitation for online alignment.
- Zhang, S., Yu, D., Sharma, H., Yang, Z., Wang, S., Hassan, H., and Wang, Z. (2024b). Self-exploring language models: Active preference elicitation for online alignment. *arXiv preprint arXiv:2405.19332*.
- Zhang, T. (2023). *Mathematical analysis of machine learning algorithms*. Cambridge University Press.
- Zhao, Y., Joshi, R., Liu, T., Khalman, M., Saleh, M., and Liu, P. J. (2023). Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*.
- Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al. (2024). Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Zhong, H., Feng, G., Xiong, W., Zhao, L., He, D., Bian, J., and Wang, L. (2024). Dpo meets ppo: Reinforced token optimization for rlhf. *arXiv preprint arXiv:2404.18922*.
- Zhou, W., Agrawal, R., Zhang, S., Indurthi, S. R., Zhao, S., Song, K., Xu, S., and Zhu, C. (2024). Wpo: Enhancing rlhf with weighted preference optimization. *arXiv preprint arXiv:2406.11827*.
- Zhu, B., Jordan, M., and Jiao, J. (2023). Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In *International Conference on Machine Learning*, pages 43037–43067. PMLR.
- Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., and Irving, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*.

## A TOKEN-LEVEL VPO

**Token-level MDP and preference modeling.** Recall that in LLMs, the prompt  $x$  can be broken into a sequence of tokens, e.g.,  $x = (x_0, \dots, x_m)$ , from a fixed discrete vocabulary  $\mathcal{A}$ . We define the token-level MDP as a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r^*, H)$ , where  $H$  is the horizon length, i.e., the longest possible number of tokens in a sentence. The state space  $\mathcal{S}$  consists of all the possible token combinations of length  $H$ , and the transition kernel is deterministically defined as follows.

1. The initial state is defined by the prompt  $x$  as  $s_0 = \{x_0, \dots, x_m\}$ ;
2. Given the response  $y = \{y_0, \dots, y_{i-1}\}$  up to the  $i$ -th token, the state at step  $i$  is defined as  $s_i = \{x_0, \dots, x_m, y_0, \dots, y_{i-1}\}$ ;
3. Upon an action of the LLM for the next token  $a_i = y_i$ , the next state at the token-level MDP deterministically becomes  $s_{i+1} = (s_i, a_i) = (x_0, \dots, x_m, y_0, \dots, y_i)$ .

We assume that the last token of a sentence, the EOS token, is absorbing, such that the token-level MDP stays in the corresponding state as soon as the last action is the EOS token. With slight abuse of notation from earlier sections, the reward function  $r^*(s, a)$  defines the ground truth reward at state  $s$  upon action  $a$ .

Given a pair of trajectories  $\tau_1 = \{s_0, a_0^1, \dots, s_{H-1}^1, a_{H-1}^1, s_H^1\}$  and  $\tau_2 = \{s_0, a_0^2, \dots, s_{H-1}^2, a_{H-1}^2, s_H^2\}$ , the corresponding Bradley-Terry preference model (Bradley and Terry, 1952) is

$$\begin{aligned} \mathbb{P}(\tau_1 \succ \tau_2) &= \frac{\exp\left(\sum_{i=0}^{H-1} r^*(s_i^1, a_i^1)\right)}{\exp\left(\sum_{i=0}^{H-1} r^*(s_i^1, a_i^1)\right) + \exp\left(\sum_{i=0}^{H-1} r^*(s_i^2, a_i^2)\right)} \\ &= \sigma\left(\sum_{i=0}^{H-1} r^*(s_i^1, a_i^1) - \sum_{i=0}^{H-1} r^*(s_i^2, a_i^2)\right), \end{aligned}$$

A preference data sample is denoted by a tuple  $(x, \tau_+, \tau_-)$ , where  $\tau_+$  (resp.  $\tau_-$ ) is the preferred (resp. unpreferred) answer in the comparison. Given a preference dataset  $\mathcal{D}$  composed of independent samples, The negative log-likelihood can be defined as

$$\ell(r, \mathcal{D}) := - \sum_{(\tau_+, \tau_-) \in \mathcal{D}} \log \sigma \left( \sum_{i=0}^{H-1} r(s_i^+, a_i^+) - \sum_{i=0}^{H-1} r(s_i^-, a_i^-) \right), \quad (22)$$

where  $(s_i^+, a_i^+)$  (resp.  $(s_i^-, a_i^-)$ ) are the state-action pairs in the trajectory  $\tau_+$  (resp.  $\tau_-$ ).

**Token-level RL fine-tuning.** Let the entropy of policy  $\pi$  under initial state distribution  $s_0 \sim \rho$  be defined as

$$\mathcal{H}(\rho, \pi) := - \mathbb{E}_{\substack{s_0 \sim \rho, \\ a_i \sim \pi(\cdot|s_i)}} \left[ \sum_{i=0}^{H-1} \log \pi(a_i|s_i) \right],$$

and  $\mathcal{H}(s, \pi)$  be the entropy when the initial state  $s_0 = s$ . Given a reward function  $r$ , we define the KL-constrained RL objective against the reference policy  $\pi_{\text{ref}}$  as (Rafailov et al., 2024):

$$\pi_r := \operatorname{argmax}_{\pi} J(r, \pi) := \mathbb{E}_{\substack{s_0 \sim \rho, \\ a_i \sim \pi(\cdot|s_i)}} \left[ \sum_{i=0}^{H-1} \underbrace{(r(s_i, a_i) + \beta \log \pi_{\text{ref}}(a_i|s_i))}_{r_{\beta}(s_i, a_i)} \right] + \beta \mathcal{H}(\rho, \pi), \quad (23)$$

where we denote  $r_{\beta} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  as follows:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad r_{\beta}(s, a) := r(s, a) + \beta \log \pi_{\text{ref}}(a|s), \quad (24)$$

which can be seen as the actual token-wise reward function optimized by the LLM.

**Token-level DPO.** The KL-constrained RL objective (23) has a closed-form solution (Nachum et al., 2017; Cen et al., 2022) given by

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad \pi_r(a|s) = \exp((Q_{\beta}^*(s, a) - V_{\beta}^*(s))/\beta), \quad (25)$$

where  $V_{\beta}^* : \mathcal{S} \rightarrow \mathbb{R}$  and  $Q_{\beta}^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  are the optimal soft value and Q functions, respectively,

$$\forall s \in \mathcal{S} : \quad V_{\beta}^*(s) := \mathbb{E}_{a_i \sim \pi_r(\cdot|s_i)} \left[ \sum_{i=0}^{H-1} r_{\beta}(s_i, a_i) | s_0 = s \right] + \beta \mathcal{H}(s, \pi_r) \quad (26)$$

denote the optimal soft value function w.r.t. the reward function  $r_{\beta}$ , and

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad Q_{\beta}^*(s, a) := r_{\beta}(s, a) + V_{\beta}^*(s'), \quad (27)$$

where  $s' = (s, a)$  is the deterministic next state. Plugging (27) and (24) into (25) implies that, for any trajectory  $\tau = \{s_0, a_0, \dots, a_{H-1}, s_H\}$ , Rafailov et al. (2024) shows

$$\begin{aligned} \sum_{i=0}^{H-1} r(s_i, a_i) &= \sum_{i=0}^{H-1} (Q_{\beta}^*(s_i, a_i) - \beta \log \pi_{\text{ref}}(a_i|s_i) - V_{\beta}^*(s_{i+1})) \\ &= Q_{\beta}^*(s_0, a_0) - \beta \log \pi_{\text{ref}}(a_0|s_0) + \sum_{i=1}^{H-1} (Q_{\beta}^*(s_i, a_i) - V_{\beta}^*(s_i) - \beta \log \pi_{\text{ref}}(a_i|s_i)) \\ &= V_{\beta}^*(s_0) + \beta \sum_{i=0}^{H-1} \log \frac{\pi_r(a_i|s_i)}{\pi_{\text{ref}}(a_i|s_i)}, \end{aligned} \quad (28)$$

where the second line uses  $V_{\beta}^*(s_H) = 0$  at the terminal state. Thus the DPO loss (which is the negative log-likelihood loss) could be written as

$$\mathcal{L}(\pi, \mathcal{D}) = - \sum_{(\tau_+, \tau_-) \in \mathcal{D}} \log \sigma \left( \beta \sum_{i=0}^{H-1} \log \frac{\pi(a_i^+|s_i^+)}{\pi_{\text{ref}}(a_i^+|s_i^+)} - \beta \sum_{i=0}^{H-1} \log \frac{\pi(a_i^-|s_i^-)}{\pi_{\text{ref}}(a_i^-|s_i^-)} \right). \quad (29)$$

**Token-level VPO.** With slight abuse of notation, define

$$J^*(r) := \max_{\pi} J(r, \pi), \quad (30)$$

which is used as the bias term in regularizing the reward estimation in VPO. Again, we impose the following assumption to deal with the shift ambiguity issue caused by the Bradley-Terry model:

**Assumption 3** We assume that  $r^* \in \mathcal{R}$ , where

$$\mathcal{R} = \left\{ r : \mathbb{E}_{\substack{s_0 \sim \rho, \\ a_i \sim \pi_{\text{cal}}(\cdot | s_i)}} \sum_{i=0}^{H-1} r(s_i, a_i) = 0 \right\}. \quad (31)$$

Here,  $\rho$  is the prompt distribution and  $\pi_{\text{cal}}$  is a fixed calibration distribution independent of the algorithm.

Combining (23) with (26), similar to previous derivations, we have

$$\begin{aligned} J^*(r) &= \mathbb{E}_{s_0 \sim \rho} [V_\beta^*(s_0)] \\ &= \mathbb{E}_{\substack{s_0 \sim \rho, \\ a_i \sim \pi_{\text{cal}}(\cdot | s_i)}} [V_\beta^*(s_0)] \\ &= \mathbb{E}_{\substack{s_0 \sim \rho, \\ a_i \sim \pi_{\text{cal}}(\cdot | s_i)}} \left[ \sum_{i=0}^{H-1} r(s_i, a_i) - \beta \sum_{i=0}^{H-1} \log \frac{\pi_r(a_i | s_i)}{\pi_{\text{ref}}(a_i | s_i)} \right] \\ &= -\beta \mathbb{E}_{\substack{s_0 \sim \rho, \\ a_i \sim \pi_{\text{cal}}(\cdot | s_i)}} \left[ \sum_{i=0}^{H-1} \log \frac{\pi_r(a_i | s_i)}{\pi_{\text{ref}}(a_i | s_i)} \right], \end{aligned} \quad (32)$$

where the penultimate line uses (28), and the last line uses Assumption 3.

Consequently, the token-level VPO can be rewritten as

$$\begin{aligned} \pi_{\text{VPO}} = \arg \min_{\pi} \left\{ - \sum_{(\tau_+, \tau_-) \in \mathcal{D}} \log \sigma \left( \beta \sum_{i=0}^{H-1} \log \frac{\pi(a_i^+ | s_i^+)}{\pi_{\text{ref}}(a_i^- | s_i^-)} - \beta \sum_{i=0}^{H-1} \log \frac{\pi(a_i^- | s_i^-)}{\pi_{\text{ref}}(a_i^- | s_i^-)} \right) \right. \\ \left. + \text{sign} \cdot \alpha \beta \mathbb{E}_{\substack{s_0 \sim \rho, \\ a_i \sim \pi_{\text{cal}}(\cdot | s_i)}} \left[ \sum_{i=0}^{H-1} \log \frac{\pi(a_i | s_i)}{\pi_{\text{ref}}(a_i | s_i)} \right] \right\}. \end{aligned} \quad (33)$$

## B ANALYSIS FOR THE ONLINE SETTING

### B.1 PROOF OF THEOREM 1

For ease of presentation, we assume that  $\mathcal{R}$  is finite, i.e.,  $|\mathcal{R}| < \infty$ . The general case can be directly obtained using a covering number argument, which we refer to (Liu et al., 2024a; Jin et al., 2022) for interested readers.

We start by decomposing the regret into two parts:

$$\begin{aligned} \text{Regret} &:= \sum_{t=1}^T [J^*(r^*) - J(r^*, \pi^{(t)})] \\ &= \underbrace{\sum_{t=1}^T [J^*(r^*) - J^*(r^{(t)})]}_{\text{Term (i)}} + \underbrace{\sum_{t=1}^T [J(r^{(t)}, \pi^{(t)}) - J(r^*, \pi^{(t)})]}_{\text{Term (ii)}}. \end{aligned} \quad (34)$$

**Step 1: bounding term (i).** By the choice of  $r^{(t)}$ , we have

$$\ell(r^{(t)}, \mathcal{D}^{(t-1)}) - \alpha J^*(r^{(t)}) \leq \ell(r^*, \mathcal{D}^{(t-1)}) - \alpha J^*(r^*). \quad (35)$$

Rearranging terms,

$$J^*(r^*) - J^*(r^{(t)}) \leq \frac{1}{\alpha} [\ell(r^*, \mathcal{D}^{(t-1)}) - \ell(r^{(t)}, \mathcal{D}^{(t-1)})]. \quad (36)$$

The following lemma is adapted from (Liu et al., 2024a, Proposition 5.3), whose proof is deferred to Appendix B.2.

**Lemma 2** Let  $\delta \in (0, 1)$ . With probability  $1 - \delta$ , we have

$$\ell(r^*, \mathcal{D}^{(t-1)}) - \ell(r^{(t)}, \mathcal{D}^{(t-1)})$$



$$\leq -2 \sum_{s=1}^{t-1} \mathbb{E}_{\substack{x \sim \rho, \\ (y_1, y_2) \sim \pi^{(s)}(\cdot|x)}} [D_H^2(\mathbb{P}_{r^{(t)}}(\cdot|x, y_1, y_2) \| \mathbb{P}_{r^*}(\cdot|x, y_1, y_2))] + 2 \log(|\mathcal{R}|/\delta). \quad (37)$$

Here,  $D_H(\cdot|\cdot)$  is the Hellinger distance,  $\mathbb{P}_r(\cdot|x, y_1, y_2)$  denotes the Bernoulli distribution of the comparison result of  $(x, y_1)$  and  $(x, y_2)$  under reward model  $r$ .

Putting the above inequalities together, it holds with probability  $1 - \delta$  that

$$\begin{aligned} \text{Term (i)} &\leq -\frac{2}{\alpha} \sum_{t=1}^T \sum_{s=1}^{t-1} \mathbb{E}_{\substack{x^{(s)} \sim \rho, \\ (y_1^{(s)}, y_2^{(s)}) \sim \pi^{(s)}(\cdot|x^{(s)})}} \left[ D_H^2(\mathbb{P}_{r^{(t)}}(\cdot|x^{(s)}, y_1^{(s)}, y_2^{(s)}) \| \mathbb{P}_{r^*}(\cdot|x^{(s)}, y_1^{(s)}, y_2^{(s)})) \right] \\ &\quad + 2\alpha^{-1} T \log(|\mathcal{R}|/\delta). \end{aligned} \quad (38)$$

**Step 2: breaking down term (ii) with the elliptical potential lemma.** The linear function approximation form (15) allows us to write

$$\mathbb{E}_{x \sim \rho, y \sim \pi_{r_\theta}(\cdot|x)} [r_1(x, y) - r^*(x, y)] = \langle W(r_\theta), X(r_\theta) \rangle, \quad (39)$$

where  $X, W : \mathcal{R} \rightarrow \mathbb{R}^d$  is given by

$$X(r_\theta) = 2C \mathbb{E}_{x \sim \rho, y \sim \pi_{r_\theta}(\cdot|x)} [\phi(x, y)], \quad W(r_\theta) = \frac{\theta - \theta^*}{2C}. \quad (40)$$

Let

$$\Sigma_t = \epsilon I + \sum_{s=1}^{t-1} X(r^{(s)}) X(r^{(s)})^\top \quad (41)$$

for some  $\epsilon > 0$ . We begin by decomposing term (ii) as

$$\begin{aligned} \text{Term (ii)} &= \sum_{t=1}^T \mathbb{E}_{x \sim \rho, y \sim \pi^{(t)}(\cdot|x)} \left[ r^{(t)}(x, y) - r^*(x, y) \right] \\ &= \sum_{t=1}^T \langle W(r^{(t)}), X(r^{(t)}) \rangle \\ &= \sum_{t=1}^T \langle W(r^{(t)}), X(r^{(t)}) \rangle \mathbf{1}\{\|X(r^{(t)})\|_{\Sigma_t^{-1}} \leq 1\} \\ &\quad + \sum_{t=1}^T \langle W(r^{(t)}), X(r^{(t)}) \rangle \mathbf{1}\{\|X(r^{(t)})\|_{\Sigma_t^{-1}} > 1\}, \end{aligned} \quad (42)$$

where  $\mathbf{1}\{A\}$  is an indicator function of event  $A$ . To proceed, we recall the elliptical potential lemma for controlling the cumulative sum of  $\min\{\|X(r^{(t)})\|_{\Sigma_t^{-1}}^2, 1\}$ .

**Lemma 3 ((Abbasi-Yadkori et al., 2011, Lemma 11))** Let  $\{X_t\}$  be a sequence in  $\mathbb{R}^d$  and  $\Lambda_0 \in \mathbb{R}^{d \times d}$  a positive definite matrix. Define  $\Lambda_t = \Lambda_0 + \sum_{s=1}^t X_s X_s^\top$ . Assume  $\|X_t\| \leq L$  for all  $t$ . It holds that

$$\begin{aligned} \sum_{t=1}^T \min\{\|X_t\|_{\Lambda_t^{-1}}^2, 1\} &\leq 2 \log \left( \frac{\det(\Lambda_T)}{\det(\Lambda_0)} \right) \\ &\leq 2(d \log((\text{trace}(\Lambda_0) + TL^2)/d) - \log \det(\Lambda_0)). \end{aligned}$$

Applying the above lemma yields

$$\sum_{t=1}^T \min\{\|X(r^{(t)})\|_{\Sigma_t^{-1}}^2, 1\} \leq \min \left\{ 2d \log \left( \frac{4C^4 T/d + \epsilon}{\epsilon} \right), T \right\} := d(\epsilon). \quad (43)$$

We now control the two terms in (42).

- The first term of (42) can be bounded by

$$\begin{aligned}
& \sum_{t=1}^T \langle W(r^{(t)}), X(r^{(t)}) \rangle \mathbf{1}\{\|X(r^{(t)})\|_{\Sigma_t^{-1}} \leq 1\} \\
& \leq \sum_{t=1}^T \|W(r^{(t)})\|_{\Sigma_t} \|X(r^{(t)})\|_{\Sigma_t^{-1}} \mathbf{1}\{\|X(r^{(t)})\|_{\Sigma_t^{-1}} \leq 1\} \\
& \leq \sum_{t=1}^T \|W(r^{(t)})\|_{\Sigma_t} \min\{\|X(r^{(t)})\|_{\Sigma_t^{-1}}, 1\} \\
& = \sum_{t=1}^T \left[ \epsilon \|W(r^{(t)})\|_2^2 + \sum_{s=1}^{t-1} \langle W(r^{(t)}), X(r^{(s)}) \rangle^2 \right]^{1/2} \min\{\|X(r^{(t)})\|_{\Sigma_t^{-1}}^2, 1\}^{1/2} \\
& \stackrel{(i)}{\leq} \left\{ \sum_{t=1}^T \left[ \epsilon \|W(r^{(t)})\|_2^2 + \sum_{s=1}^{t-1} \langle W(r^{(t)}), X(r^{(s)}) \rangle^2 \right] \right\}^{1/2} \left\{ \sum_{t=1}^T \min\{\|X(r^{(t)})\|_{\Sigma_t^{-1}}^2, 1\} \right\}^{1/2} \\
& \stackrel{(ii)}{\leq} \sqrt{d(\epsilon)} \left\{ \sum_{t=1}^T \sum_{s=1}^{t-1} \langle W(r^{(t)}), X(r^{(s)}) \rangle^2 \right\}^{1/2} + \sqrt{d(\epsilon)\epsilon T} \\
& \stackrel{(iii)}{\leq} \frac{d(\epsilon)}{2\mu} + \frac{\mu}{2} \sum_{t=1}^T \sum_{s=1}^{t-1} \langle W(r^{(t)}), X(r^{(s)}) \rangle^2 + \sqrt{d(\epsilon)\epsilon T}. \tag{44}
\end{aligned}$$

Here, (i) is due to Cauchy–Schwarz inequality, (ii) is due to  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for  $\forall a, b \geq 0$ , and (iii) results from Young’s inequality. We leave the constant  $\mu > 0$  to be determined later.

- The second term of (42) can be bounded by

$$\begin{aligned}
& \sum_{t=1}^T \langle W(r^{(t)}), X(r^{(t)}) \rangle \mathbf{1}\{\|X(r^{(t)})\|_{\Sigma_t^{-1}} > 1\} \leq C \sum_{t=1}^T \mathbf{1}\{\|X(r^{(t)})\|_{\Sigma_t^{-1}} > 1\} \\
& \leq C \sum_{t=1}^T \min\{\|X(r^{(t)})\|_{\Sigma_t^{-1}}^2, 1\} \leq Cd(\epsilon), \tag{45}
\end{aligned}$$

where the first inequality follows from  $\|X(r^{(t)})\|_2 \leq 2C$  and  $\|W(r^{(t)})\|_2 \leq 1/2$  since  $\|\phi(x, y)\|_2 \leq 1$ .

Putting (42), (44) and (45) together, we arrive at

$$\text{Term (ii)} \leq \frac{d(\epsilon)}{2\mu} + \frac{\mu}{2} \sum_{t=1}^T \sum_{s=1}^{t-1} \langle W(r^{(t)}), X(r^{(s)}) \rangle^2 + \sqrt{d(\epsilon)\epsilon T} + Cd(\epsilon). \tag{46}$$

**Step 3: continuing bounding term (ii).** It boils down to control  $\langle W(r^{(t)}), X(r^{(s)}) \rangle^2$ . We have

$$\begin{aligned}
\langle W(r^{(t)}), X(r^{(s)}) \rangle & = \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^{(s)}(\cdot|x)}} \left[ r^{(t)}(x, y) - r^*(x, y) \right] \\
& = \mathbb{E}_{\substack{x \sim \rho, \\ y_1 \sim \pi^{(s)}(\cdot|x)}} \left[ r^{(t)}(x, y_1) - r^*(x, y_1) \right] - \mathbb{E}_{\substack{x \sim \rho, \\ y_2 \sim \pi_{\text{cal}}(\cdot|x)}} \left[ r^{(t)}(x, y_2) - r^*(x, y_2) \right] \\
& = \mathbb{E}_{\substack{x \sim \rho, \\ y_1 \sim \pi^{(s)}(\cdot|x), \\ y_2 \sim \pi_{\text{cal}}(\cdot|x)}} \left[ \delta_x(r^{(t)}, r^*, y_1, y_2) \right], \tag{47}
\end{aligned}$$

where  $\delta_x(r_1, r_2, y_1, y_2) := r_1(x, y_1) - r_1(x, y_2) - (r_2(x, y_1) - r_2(x, y_2))$ . Therefore,

$$\langle W(r^{(t)}), X(r^{(s)}) \rangle^2 = \mathbb{E}_{\substack{x \sim \rho, \\ y_1 \sim \pi^{(s)}(\cdot|x), \\ y_2 \sim \pi_{\text{cal}}(\cdot|x)}} \left[ \delta_x(r^{(t)}, r^*, y_1, y_2) \right]^2$$

$$\begin{aligned}
&= \mathbb{E}_{\substack{x \sim \rho, \\ y_1 \sim \pi^{(s)}(\cdot|x), \\ y_2 \sim \pi_{\text{cal}}(\cdot|x)}} \left[ \delta_x(r^{(t)}, r^*, y_1, y_2)^2 \right] - \text{Var}_{\substack{x \sim \rho, \\ y_1 \sim \pi^{(s)}(\cdot|x), \\ y_2 \sim \pi_{\text{cal}}(\cdot|x)}} \left[ \delta_x(r^{(t)}, r^*, y_1, y_2)^2 \right] \\
&\leq \mathbb{E}_{\substack{x \sim \rho, \\ y_1 \sim \pi^{(s)}(\cdot|x), \\ y_2 \sim \pi_{\text{cal}}(\cdot|x)}} \left[ \delta_x(r^{(t)}, r^*, y_1, y_2)^2 \right] \\
&\leq \sup_{x,y} \frac{\pi_{\text{cal}}(y|x)}{\pi^{(s)}(y|x)} \cdot \mathbb{E}_{\substack{x \sim \rho, \\ y_1, y_2 \sim \pi^{(s)}(\cdot|x)}} \left[ \delta_x(r^{(t)}, r^*, y_1, y_2)^2 \right] \\
&\leq \sup_{x,y} \frac{\pi_{\text{ref}}(y|x)}{\pi^{(s)}(y|x)} \cdot \sup_{x,y} \frac{\pi_{\text{cal}}(y|x)}{\pi_{\text{ref}}(y|x)} \cdot \mathbb{E}_{\substack{x \sim \rho, \\ y_1, y_2 \sim \pi^{(s)}(\cdot|x)}} \left[ \delta_x(r^{(t)}, r^*, y_1, y_2)^2 \right]. \quad (48)
\end{aligned}$$

Recall from (4) that  $\pi^{(s)}(y|x) \propto \pi_{\text{ref}}(y|x) \exp(r^{(s)}(x, y)/\beta)$ . It follows that  $|\log \pi^{(s)}(y|x) - \log \pi_{\text{ref}}(y|x)| \leq 2\|r^{(s)}(x, \cdot)\|_\infty \leq 2C/\beta$  (see e.g., (Cen et al., 2022, Appendix A.2)), and hence  $\sup_{x,y} \frac{\pi_{\text{ref}}(y|x)}{\pi^{(s)}(y|x)} \leq \exp(2C/\beta)$ . To proceed, we demonstrate in the following lemma that  $\delta^2$  can be upper bounded by the corresponding Hellinger distance, whose proof is deferred to Appendix B.3.

**Lemma 4** Assume bounded reward  $\|r_1\|_\infty \leq C, \|r_2\|_\infty \leq C$ . We have

$$\delta_x(r_1, r_2, y_1, y_2)^2 \leq 2(3 + \exp(2C))^2 D_{\text{H}}^2(\mathbb{P}_{r_1}(\cdot|x, y_1, y_2) \parallel \mathbb{P}_{r_2}(\cdot|x, y_1, y_2)).$$

With the above lemma we arrive at

$$\begin{aligned}
&\langle W(r^{(t)}), X(r^{(s)}) \rangle^2 \\
&\leq 2(3 + \exp(2C))^2 \exp(2C/\beta) \kappa \cdot \mathbb{E}_{\substack{x \sim \rho, \\ y_1, y_2 \sim \pi^{(s)}(\cdot|x)}} \left[ D_{\text{H}}^2(\mathbb{P}_{r^{(t)}}(\cdot|x, y_1, y_2) \parallel \mathbb{P}_{r^*}(\cdot|x, y_1, y_2)) \right].
\end{aligned}$$

where we denote  $\kappa = \sup_{x,y} \frac{\pi_{\text{cal}}(y|x)}{\pi_{\text{ref}}(y|x)}$ . Plugging the above bound into (46), we get

**Term (ii)**

$$\begin{aligned}
&\leq \frac{d(\epsilon)}{2\mu} + \mu(3 + \exp(2C))^2 \exp(2C/\beta) \kappa \cdot \sum_{t=1}^T \sum_{s=1}^{t-1} \mathbb{E}_{\substack{x \sim \rho, \\ y_1, y_2 \sim \pi^{(s)}(\cdot|x)}} \left[ D_{\text{H}}^2(\mathbb{P}_{r^{(t)}}(\cdot|x, y_1, y_2) \parallel \mathbb{P}_{r^*}(\cdot|x, y_1, y_2)) \right] \\
&\quad + 2B\sqrt{d(\epsilon)\epsilon T} + Cd(\epsilon). \quad (49)
\end{aligned}$$

**Step 4: finishing up.** Combining (34), (38) and (49), with probability  $1 - \delta$  we have

$$\text{Regret} \leq \frac{2T \log(|\mathcal{R}|/\delta)}{\alpha} + \frac{d(\epsilon)}{2\mu} + \sqrt{d(\epsilon)\epsilon T} + Cd(\epsilon) \quad (50)$$

as long as  $\alpha\mu(3 + \exp(2C))^2 \exp(2C/\beta) \kappa \leq 2$ . Setting  $\alpha \asymp \frac{1}{\exp(2C+C/\beta)} \sqrt{\frac{T}{\kappa d(\epsilon)}}$ ,  $\mu \asymp \frac{1}{\exp(2C+C/\beta)} \sqrt{\frac{d(\epsilon)}{\kappa T}}$ , and  $\epsilon = 1$ , we arrive at

$$\text{Regret} \leq \tilde{O}((\exp(2C + C/\beta))\sqrt{\kappa d T})$$

as claimed.

## B.2 PROOF OF LEMMA 2

To begin, we have

$$\ell(r^*, \mathcal{D}^{(t-1)}) - \ell(r^{(t)}, \mathcal{D}^{(t-1)}) = -\log \frac{\mathbb{P}(\mathcal{D}^{(t-1)}|r^*)}{\mathbb{P}(\mathcal{D}^{(t-1)}|r^{(t)})} = -\sum_{s=1}^{t-1} X_r^s, \quad (51)$$

where we denote

$$X_r^s = \log \frac{\mathbb{P}_{r^*}(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})}{\mathbb{P}_r(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})}. \quad (52)$$

To proceed, we recall a useful martingale exponential inequality.

**Lemma 5** ((Zhang, 2023, Theorem 13.2), (Liu et al., 2024a, Lemma D.1)) *Let  $\{X_t\}_{t=1}^\infty$  be a sequence of real-valued random variables adapted to filtration  $\{\mathcal{F}_t\}_{t=1}^\infty$ . It holds with probability  $1 - \delta$  such that for any  $t \geq 1$ ,*

$$-\sum_{s=1}^t X_s \leq \sum_{s=1}^t \log \mathbb{E} [\exp(-X_s) | \mathcal{F}_{s-1}] + \log(1/\delta).$$

Applying the above lemma to  $\{\frac{1}{2}X_r^s\}_{s=1}^\infty$  along with the filtration  $\{\mathcal{F}_t\}_{t=1}^\infty$  with  $\mathcal{F}_t$  given by the  $\sigma$ -algebra of  $\{(x^{(s)}, y_+^{(s)}, y_-^{(s)}) : s \leq t\}$ , we conclude that it holds with probability  $1 - \delta$  that

$$\begin{aligned} -\frac{1}{2} \sum_{s=1}^{t-1} X_r^s &\leq \sum_{s=1}^{t-1} \log \mathbb{E} \left[ \exp \left\{ -\frac{1}{2} X_r^s \right\} \middle| \mathcal{F}_{s-1} \right] + \log(|\mathcal{R}|/\delta) \\ &\leq \sum_{s=1}^{t-1} \left( \mathbb{E} \left[ \exp \left\{ -\frac{1}{2} X_r^s \right\} \middle| \mathcal{F}_{s-1} \right] - 1 \right) + \log(|\mathcal{R}|/\delta), \end{aligned} \quad (53)$$

where the last step results from the inequality  $\log(1+x) \leq x$  for all  $x \geq -1$ . To proceed, note that

$$\begin{aligned} &\mathbb{E} \left[ \exp \left\{ -\frac{1}{2} X_r^s \right\} \middle| \mathcal{F}_{s-1} \right] \\ &= \mathbb{E} \left[ \sqrt{\frac{\mathbb{P}_r(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})}{\mathbb{P}_{r^*}(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})}} \middle| \mathcal{F}_{s-1} \right] \\ &= \mathbb{E}_{\substack{x^{(s)} \sim \rho, \\ (y_1^{(s)}, y_2^{(s)}) \sim \pi^{(s)}(\cdot | x^{(s)}), \\ (+, -) \sim \mathbb{P}_{r^*}}} \left[ \sqrt{\frac{\mathbb{P}_r(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})}{\mathbb{P}_{r^*}(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})}} \right] \\ &= \mathbb{E}_{\substack{x^{(s)} \sim \rho, \\ (y_1^{(s)}, y_2^{(s)}) \sim \pi^{(s)}(\cdot | x^{(s)})}} \left[ \sum_{(+, -)} \sqrt{\mathbb{P}_r(y_+^{(s)} \succ y_-^{(s)} | x^{(s)}) \cdot \mathbb{P}_{r^*}(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})} \right] \\ &= 1 - \frac{1}{2} \mathbb{E}_{\substack{x^{(s)} \sim \rho, \\ (y_1^{(s)}, y_2^{(s)}) \sim \pi^{(s)}(\cdot | x^{(s)})}} \left[ \sum_{(+, -)} \left( \sqrt{\mathbb{P}_r(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})} - \sqrt{\mathbb{P}_{r^*}(y_+^{(s)} \succ y_-^{(s)} | x^{(s)})} \right)^2 \right] \\ &= 1 - \mathbb{E}_{\substack{x \sim \rho, \\ (y_1, y_2) \sim \pi^{(s)}(\cdot | x)}} \left[ D_{\text{H}}^2(\mathbb{P}_r(\cdot | x, y_1, y_2) \| \mathbb{P}_{r^*}(\cdot | x, y_1, y_2)) \right], \end{aligned}$$

where we denote by  $\sum_{(+, -)}$  the summation over different comparison results. Plugging the above equality into (53) completes the proof.

### B.3 PROOF OF LEMMA 4

By the mean value theorem, we have

$$\begin{aligned} |\mathbb{P}_{r_1}(y_1 \succ y_2 | x) - \mathbb{P}_{r_2}(y_1 \succ y_2 | x)| &= |\sigma(r_1(x, y_1) - r_1(x, y_2)) - \sigma(r_2(x, y_1) - r_2(x, y_2))| \\ &= |\delta_x(r_1, r_2, y_1, y_2) \cdot \sigma'(\xi)| \\ &= |\delta_x(r_1, r_2, y_1, y_2)| \cdot \sigma(\xi)(1 - \sigma(\xi)) \end{aligned}$$

for some  $\xi$  between  $r_1(x, y_1) - r_1(x, y_2)$  and  $r_2(x, y_1) - r_2(x, y_2)$ . Since  $|\xi| \leq 2C$ , we have

$$\sigma(\xi)(1 - \sigma(\xi)) \geq \sigma(2C)(1 - \sigma(2C)) \geq \frac{1}{3 + \exp(2C)}. \quad (54)$$

Putting together,

$$|\delta_x(r_1, r_2, y_1, y_2)| \leq (3 + \exp(2C)) |\mathbb{P}_{r_1}(y_1 \succ y_2 | x) - \mathbb{P}_{r_2}(y_1 \succ y_2 | x)|$$

$$\begin{aligned}
&= (3 + \exp(2C))\text{TV}(\mathbb{P}_{r_1}(\cdot|x, y_1, y_2), \mathbb{P}_{r_2}(\cdot|x, y_1, y_2)) \\
&\leq (3 + \exp(2C))\sqrt{2}D_{\text{H}}(\mathbb{P}_{r_1}(\cdot|x, y_1, y_2) \parallel \mathbb{P}_{r_2}(\cdot|x, y_1, y_2)).
\end{aligned}$$

## C ANALYSIS FOR THE OFFLINE SETTING

### C.1 PROOF OF LEMMA 1

By definition, the objective function  $\ell(r, \mathcal{D}) + \alpha J(r, \pi)$  is strongly concave over  $\pi$ , and convex over  $r$ . By Danskin's theorem, we have

$$\nabla_r(\max_{\pi}[\ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \pi)]) = \nabla_r(\ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi})).$$

Therefore, for any  $r'$ , by convexity of the objective function we have

$$\begin{aligned}
\ell(r', \mathcal{D}) + \alpha J(r', \hat{\pi}) &\geq \ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi}) + \langle r' - \hat{r}, \nabla_r(\ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi})) \rangle \\
&= \ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi}) + \langle r' - \hat{r}, \nabla_r(\max_{\pi}[\ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \pi)]) \rangle \\
&\geq \ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi}).
\end{aligned}$$

The last line is due to the definition of  $\hat{r}$  (c.f. (20)). The other relation,  $\ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \hat{\pi}) \geq \ell(\hat{r}, \mathcal{D}) + \alpha J(\hat{r}, \pi')$ , follows directly from the definition of  $\hat{\pi}$  (c.f. (17)).

### C.2 PROOF OF THEOREM 2

We decompose the sub-optimality gap of  $\hat{\pi}$  by

$$\begin{aligned}
&J^*(r^*) - J(r^*, \hat{\pi}) \\
&= [J(r^*, \pi^*) - J(\hat{r}, \pi^*)] + [J(\hat{r}, \pi^*) - J(\hat{r}, \hat{\pi})] + [J(\hat{r}, \hat{\pi}) - J(r^*, \hat{\pi})] \\
&\leq \underbrace{[J(r^*, \pi^*) - J(\hat{r}, \pi^*)]}_{\text{Term (i)}} + \underbrace{[J(\hat{r}, \hat{\pi}) - J(r^*, \hat{\pi})]}_{\text{Term (ii)}}, \tag{55}
\end{aligned}$$

where the last line is due to  $J(\hat{r}, \pi^*) \leq J(\hat{r}, \hat{\pi})$  according to the definition of  $\hat{\pi}$  (c.f. (17)). We proceed to bound the two terms separately. Here we have written  $\hat{r} = r_{\hat{\theta}}$  for notational simplicity. In addition, we denote the MLE estimate by  $r_{\text{MLE}} = r_{\theta_{\text{MLE}}}$ .

By the definition of  $J(r, \pi)$  (cf. (3)), it follows that term (i) in (55) can be further decomposed as

$$\begin{aligned}
\text{Term (i)} &= \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [r^*(x, y) - \hat{r}(x, y)] \\
&= \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\langle \phi(x, y), \theta^* - \hat{\theta} \rangle] \\
&= \underbrace{\mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\langle \phi(x, y), \theta^* - \theta_{\text{MLE}} \rangle]}_{\text{Term (ia)}} + \underbrace{\mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\langle \phi(x, y), \theta_{\text{MLE}} - \hat{\theta} \rangle]}_{\text{Term (ib)}}, \tag{56}
\end{aligned}$$

where  $r_{\text{MLE}}(x, y) = \langle \phi(x, y), \theta_{\text{MLE}} \rangle$ .

**Step 1: bounding term (ia).** To continue, we recall a useful lemma from (Zhu et al., 2023).

**Lemma 6 ((Zhu et al., 2023, Lemma 3.1))** For any  $\lambda > 0$  and  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\|\theta_{\text{MLE}} - \theta^*\|_{\Sigma_{\mathcal{D}} + \lambda I} \leq \mathcal{O}\left(\left(3 + \exp(C)\right)\sqrt{\frac{d + \log(1/\delta)}{N}} + \sqrt{\lambda C^2}\right).$$

In addition, we have

$$\frac{1}{3 + \exp(C)}\Sigma_{\mathcal{D}} \preceq \frac{1}{N}\nabla_{\theta}^2\ell(r_{\theta}, \mathcal{D}) \preceq \frac{1}{4}\Sigma_{\mathcal{D}} \tag{57}$$

for all  $\theta$  such that  $\|r_{\theta}\|_{\infty} \leq C$ .

1134 The first term of (56) can be bounded with Lemma 6 as

$$\begin{aligned}
1135 \text{Term (ia)} &\leq \|\theta^* - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I} \cdot \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} \\
1136 &\leq \mathcal{O} \left( \left( (3 + \exp(C)) \sqrt{\frac{d + \log(1/\delta)}{N}} + \sqrt{\lambda C^2} \right) \cdot \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} \right). \\
1137 & \\
1138 & \\
1139 & \\
1140 & \\
1141 & \tag{58}
\end{aligned}$$

1142 **Step 2: bounding term (ib).** For the second term of (56), recall that

$$1143 \hat{r} = \arg \min_{r \in \mathcal{R}} \{ \ell(r, \mathcal{D}) + \alpha J(r, \hat{\pi}) \},$$

1144 or equivalently

$$1145 \hat{\theta} = \arg \min_{\theta \in \Theta} \{ \ell(r_{\theta}, \mathcal{D}) + \alpha J(r_{\theta}, \hat{\pi}) \},$$

1146 and that

$$1147 \theta_{\text{MLE}} = \arg \min_{\theta \in \Theta} \ell(r_{\theta}, \mathcal{D}).$$

1148 With linear constraint (16), by KKT condition we have

$$1149 \nabla_{\theta} \ell(\hat{r}, \mathcal{D}) + \alpha \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] + \lambda_1 \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)] = 0$$

1150 for some  $\lambda_1 \in \mathbb{R}$ , and

$$1151 \nabla_{\theta} \ell(r_{\text{MLE}}, \mathcal{D}) + \lambda_2 \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)] = 0$$

1152 for some  $\lambda_2 \in \mathbb{R}$ . By strong monotonicity of  $\nabla_{\theta} \ell$  (cf. (57)), we have

$$\begin{aligned}
1153 \frac{N}{3 + \exp(C)} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}}}^2 &\leq \langle \nabla_{\theta} \ell(\hat{r}, \mathcal{D}) - \nabla_{\theta} \ell(r_{\text{MLE}}, \mathcal{D}), \hat{\theta} - \theta_{\text{MLE}} \rangle \\
1154 &= \left\langle -\alpha \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - (\lambda_1 - \lambda_2) \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)], \hat{\theta} - \theta_{\text{MLE}} \right\rangle \\
1155 &= -\alpha \left\langle \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)], \hat{\theta} - \theta_{\text{MLE}} \right\rangle \\
1156 &\leq \alpha \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I} \\
1157 &\leq \alpha \kappa_{\mathcal{D}} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I}, \\
1158 &
\end{aligned}$$

1159 where we denote

$$1160 \kappa_{\mathcal{D}} = \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}}. \tag{59}$$

1161 The penultimate step results from  $\hat{\theta}, \theta_{\text{MLE}} \in \Theta$ , which ensures

$$1162 \left\langle \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)], \hat{\theta} \right\rangle = \left\langle \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)], \theta_{\text{MLE}} \right\rangle = 0$$

1163 It follows that

$$\begin{aligned}
1164 \frac{N}{3 + \exp(C)} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I}^2 &\leq \frac{N}{3 + \exp(C)} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}}}^2 + \frac{N}{3 + \exp(C)} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\lambda I}^2 \\
1165 &\leq \alpha \kappa_{\mathcal{D}} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I} + \frac{N \lambda C^2}{3 + \exp(C)}. \\
1166 &
\end{aligned}$$

1167 The above inequality allows us to bound

$$1168 \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I} \leq \frac{\alpha(3 + \exp(C))}{N} \kappa_{\mathcal{D}} + 2\sqrt{\lambda C^2}. \tag{60}$$

1169

1188 Therefore, the second term of (56) can be bounded as

$$\begin{aligned}
 1189 \quad \text{Term (ib)} &\leq \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I} \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} \\
 1192 &\leq \left( \frac{\alpha(3 + \exp(C))}{N} \kappa_{\mathcal{D}} + 2\sqrt{\lambda C^2} \right) \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}}. \quad (61)
 \end{aligned}$$

1193 Putting (58) and (61) together, we have

$$\begin{aligned}
 1196 \quad \text{Term (i)} &\leq \mathcal{O} \left( \left[ \frac{3 + \exp(C)}{\sqrt{N}} \left( \sqrt{d + \log(1/\delta)} + \frac{\alpha}{\sqrt{N}} \kappa_{\mathcal{D}} \right) + \sqrt{\lambda C^2} \right] \cdot \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} \right). \\
 1199 &\quad (62)
 \end{aligned}$$

1200 **Step 3: bounding term (ii).** We can decompose and bound term (ii) by

$$\begin{aligned}
 1202 \quad J(\hat{r}, \hat{\pi}) - J(r^*, \hat{\pi}) &= J(\hat{r}, \hat{\pi}) + \frac{1}{\alpha} \ell(\hat{r}, \mathcal{D}) - \left( J(r^*, \hat{\pi}) + \frac{1}{\alpha} \ell(r^*, \mathcal{D}) \right) + \frac{1}{\alpha} (\ell(\hat{r}, \mathcal{D}) - \ell(r^*, \mathcal{D})) \\
 1203 &\stackrel{(i)}{\leq} \frac{1}{\alpha} (\ell(\hat{r}, \mathcal{D}) - \ell(r^*, \mathcal{D})) \\
 1205 &\leq \frac{1}{\alpha} (\ell(\hat{r}, \mathcal{D}) - \ell(r_{\text{MLE}}, \mathcal{D}) + \ell(r_{\text{MLE}}, \mathcal{D}) - \ell(r^*, \mathcal{D})),
 \end{aligned}$$

1206 where (i) follows from the fact that  $(\hat{r}, \hat{\pi})$  is a saddle point. Due to convexity of  $\ell$ , we have

$$\begin{aligned}
 1209 \quad \ell(\hat{r}, \mathcal{D}) - \ell(r_{\text{MLE}}, \mathcal{D}) &\leq \langle \nabla_{\theta} \ell(\hat{r}, \mathcal{D}), \hat{\theta} - \theta_{\text{MLE}} \rangle \\
 1210 &= \left\langle -\alpha \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \lambda_1 \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)], \hat{\theta} - \theta_{\text{MLE}} \right\rangle \\
 1211 &= -\alpha \left\langle \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)], \hat{\theta} - \theta_{\text{MLE}} \right\rangle \\
 1212 &\leq \alpha \kappa_{\mathcal{D}} \|\hat{\theta} - \theta_{\text{MLE}}\|_{\Sigma_{\mathcal{D}} + \lambda I} \\
 1213 &\leq \frac{\alpha^2(3 + \exp(C))}{N} \kappa_{\mathcal{D}}^2 + 2\sqrt{\lambda C^2} \alpha \kappa_{\mathcal{D}},
 \end{aligned}$$

1214 where the last step is due to (60). On the other hand, with probability  $1 - \delta$  we have (Zhan et al., 2023, Lemma 1):

$$\ell(r_{\text{MLE}}, \mathcal{D}) - \ell(r^*, \mathcal{D}) \leq \tilde{\mathcal{O}}(1).$$

1221 Putting pieces together,

$$\begin{aligned}
 1222 \quad \text{Term (ii)} &\leq \frac{\alpha(3 + \exp(C))}{N} \kappa_{\mathcal{D}}^2 + 2\sqrt{\lambda C^2} \kappa_{\mathcal{D}} + \frac{1}{\alpha}. \quad (63)
 \end{aligned}$$

1223 **Step 4: putting things together.** Combining (55) (62), (63), with probability  $1 - \delta$  we have

$$\begin{aligned}
 1227 \quad J^*(r^*) - J(r^*, \hat{\pi}) &\leq \mathcal{O} \left( \frac{1}{\sqrt{N}} \left[ (3 + \exp(C)) \left( \sqrt{d + \log(1/\delta)} + \kappa_{\mathcal{D}} \right) + C \right] \cdot \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi^*(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}} \right. \\
 1229 &\quad \left. + \frac{1}{\sqrt{N}} \left( (3 + \exp(C)) \kappa_{\mathcal{D}}^2 + 2C \kappa_{\mathcal{D}} + 1 \right) \right).
 \end{aligned}$$

1232 Here we have set  $\alpha = \sqrt{N}$  and  $\lambda = 1/N$ . We conclude by bounding  $\kappa_{\mathcal{D}}$  as

$$\begin{aligned}
 1236 \quad \kappa_{\mathcal{D}}^2 &= \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)] \right\|_{(\Sigma_{\mathcal{D}} + \lambda I)^{-1}}^2 \\
 1237 &\leq \left\| \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \hat{\pi}(\cdot|x)}} [\phi(x, y)] - \mathbb{E}_{\substack{x \sim \rho, \\ y \sim \pi_{\text{cal}}(\cdot|x)}} [\phi(x, y)] \right\|_2^2 \cdot \left\| (\Sigma_{\mathcal{D}} + \lambda I)^{-1} \right\|_2 \\
 1239 &\leq 4(\lambda_{\min}(\Sigma_{\mathcal{D}}) + \lambda)^{-1}.
 \end{aligned}$$

1241

## D EXPERIMENTAL DETAILS

### D.1 RLHF FOR LLMS — OFFLINE SETTING

For the offline setting experiments, we adopt instruction tuned models, LLAMA2-7B-CHAT, LLAMA2-13B-CHAT and FLAN-T5-XL as the base models. To prompt these models, we prepend the questions in ARC-Challenge task (Clark et al., 2018) with

*What is the choice to the following Question? Only provide the choice by providing a single letter.*

and further append the question with

*The answer is:.*

The question is structured in a way that the multiple choices are appended with alphabets (letters) within parenthesis to the question. As an example:

*Question: George wants to warm his hands quickly by rubbing them. Which skin surface will produce the most heat? Choices: (A)dry palms (B)wet palms (C)palms covered with oil (D)palms covered with lotion*

We set  $\pi_{\text{cal}}$  to the empirical distribution of the ground truth answer which is known to us. Based on preliminary experiments, we set  $\beta$  as 0.1 in DPO and  $\tau$  as 1.0 in IPO. For VPO, we experiment with moving  $\alpha$  from 0.01 to 10, choosing 1 for the reported results. For all models, we train the base models with different algorithms DPO, VPO and IPO for 3000 steps and report the accuracy of the performance on the ARC-challenge test data set after every 500 steps. The training for LLAMA2-13B-CHAT model on 128 TPU-v4 takes around 2hrs and for FLAN-T5-XL on 64 TPU-v3 takes 1 hour.

### D.2 RLHF FOR LLMS — ONLINE SETTING

#### D.2.1 BUFFER

The prompt used for generating AI feedback (and rating) for TL;DR summarization is identical to (Guo et al., 2024). We set  $\pi_{\text{cal}}$  to the empirical distribution of the negative answer pairs  $(x, y_-)$  collected by the policy. We set  $\beta$  as 0.1 for the DPO term similar to (Guo et al., 2024). Additionally for VPO, we decrease the coefficient exponentially following  $\frac{\alpha}{\sqrt{1+\text{training steps}}}$ . We try different values of  $\alpha$  and report the results for 0.1 and 0.01.

The training of the policy, PALM2-XXS on 64 TPU-v3 for 5000 steps takes around 12 hours for both online DPO and VPO. We report the win rate percentage against the base SFT model for every 1000 steps using PALM2-XS judge. We also further conduct side by side comparison of Online DPO and VPO at 5000 step.

#### D.2.2 ITERATIVE

The UltraFeedback data (Cui et al., 2023) contains around 61k preference pairs of single-turn conversations. We divide this data set into 4 chunks. We use the first chunk to train a DPO model which we refer to as Zephyr-7B-DPO in 1. We use the remaining 3 chunks to train consecutive iterations of DPO and VPO, using the checkpoint from the previous iteration to initialize the policy for the current iteration. For each iteration, the prompts from the data are extracted and a new answer is sampled from the policy. We label the data using online AI feedback provided by PairRM (Jiang et al., 2023), using the similar ranking procedure as (Zhang et al., 2024b) where the new sample is ranked against  $y_w$  and  $y_l$  from the data. For VPO, we also adopt another ranking process where we sample two answer from the policy which are ranked against each other. We also set  $\pi_{\text{cal}}$  to the empirical distribution of the negative answer pairs  $(x, y_-)$  collected by the policy. We report the results of the best checkpoint of these two ranking procedures in 1. All experiments are conducted on 16xA100 GPUs.



1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

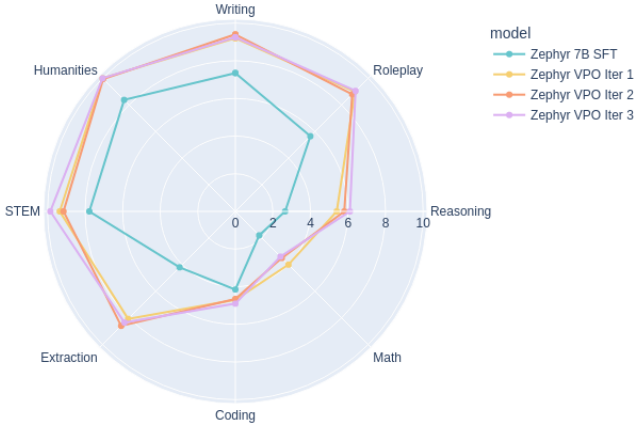


Figure 3: Radar chart of MT-Bench results for Zephyr-7B.

### D.3 SYNTHETIC BANDIT PROBLEMS

We evaluate the proposed methods on two synthetic problems: i) a multi-armed bandit problem, and ii) a linear contextual bandit problem.

**Multi-armed bandit (MAB) problem.** In this scenario, we set  $|\mathcal{X}| = 1$  and  $|\mathcal{Y}| = 10$ . For each  $y \in \mathcal{Y}$ , the ground truth reward  $r^*(x, y)$  is randomly generated i.i.d. from a uniform distribution  $U([0, 1])$ . The policy is parameterized by  $\pi_\theta(\cdot|x) = \text{softmax}(\theta(x, \cdot))$ , where  $\theta \in \mathbb{R}^{10}$ . The reference policy  $\pi_{\text{ref}}$  is set to  $\pi_{\theta_{\text{ref}}}$  with  $\theta_{\text{ref}}(x, y)$  sampled i.i.d. from  $U([0, 1])$ .

**Linear contextual bandit problem.** Here, we set  $\mathcal{X} = \mathbb{R}^2$  and  $|\mathcal{Y}| = 50$ . For each  $(x, y)$  pair, the ground truth reward is given by  $r^*(x, y) = \langle \phi(x, y), \theta^* \rangle$ , where  $\theta^* \in \mathbb{R}^{10}$  is randomly sampled from  $U([0, 1])$ , and the feature vector  $\phi(x, y)$  is the output of the hidden layer of a fixed two-layer MLP, with the input given by the concatenation of  $x$  and the one-hot encoding of  $y$ . The activation function is set to tanh. The context vector  $x$  is drawn from standard normal distribution. We focus on log-linear policy class  $\pi_\theta(\cdot|x) = \text{softmax}(\langle \theta, \phi(x, \cdot) \rangle)$ , and set  $\pi_{\text{ref}} = \pi_{\theta_{\text{ref}}}$  with  $\theta_{\text{ref}}(x, y)$  sampled i.i.d. from  $U([0, 1])$ .

For both problem we set  $\pi_{\text{ref}} = \pi_b = \pi_{\text{cal}}$  and use mini-batch sample of size 5 in every iteration. We approximately solve the optimization problems by performing 20 AdamW optimization steps with learning rate 0.01 and weight decay rate 0.01 in every iteration for the online setting and 1000 steps for the offline setting.

We plot the average results over 20 independent runs for both experiments in Figure 4 and Figure 5, as well as  $\pm 1$  standard error bars. As demonstrated in the left panel of Figure 4, an appropriate choice of  $\alpha$  allows our method to outperform the model-based MAB with MLE baseline in the long-term performance of cumulative regret, at the cost of slightly increased cumulative regret in the first 100 iterations. This highlights the effectiveness of the VPO in achieving more principled exploration-exploitation trade-off. In the right panel, it is evident that the MLE baseline struggles with principled exploration in the more complex linear contextual bandit problem, while our method is able to achieve sub-linear regret growth. For the offline setting, Figure 5 demonstrates that the performance of both MLE-MAB and VPO improves as the number of offline data increases. However, VPO achieves a consistently lower sub-optimality gap compared with that of MLE-MAB.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

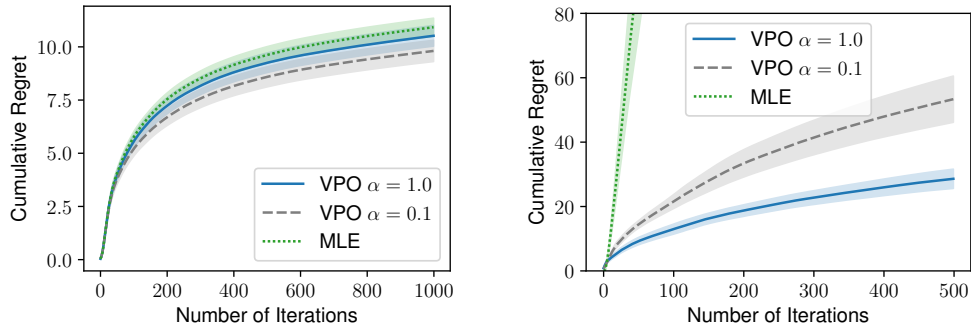


Figure 4: The cumulative regret v.s. number of iterations plot of VPO and MLE-MAB methods in the online MAB problem (left panel) and online linear contextual bandit problem (right panel), respectively.

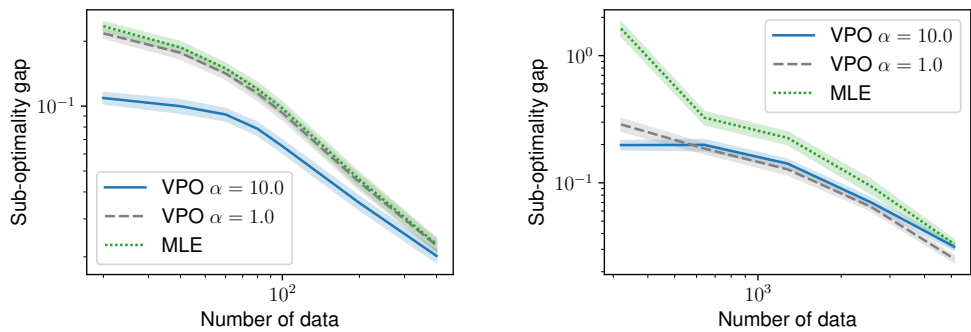


Figure 5: The sub-optimality gap v.s. number of data plot of VPO and MLE-MAB methods in the offline MAB problem (left panel) and offline linear contextual bandit problem (right panel), respectively.