

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 TOWARDS EFFICIENT ONLINE EXPLORATION FOR REINFORCEMENT LEARNING WITH HUMAN FEEDBACK

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement learning with human feedback (RLHF), which learns a reward model from human preference data and then optimizes a policy to favor preferred responses, has emerged as a central paradigm for aligning large language models (LLMs) with human preferences. In this paper, we investigate exploration principles for online RLHF, where one seeks to adaptively collect new preference data to refine both the reward model and the policy in a data-efficient manner. By examining existing optimism-based exploration algorithms, we identify a drawback in their sampling protocol: they tend to gather comparisons that fail to reduce the most informative uncertainties in reward differences, and we prove lower bounds showing that such methods can incur linear regret over exponentially long horizons. Motivated by this insight, we propose a new exploration scheme that directs preference queries toward reducing uncertainty in reward differences most relevant to policy improvement. Under a multi-armed bandit model of RLHF, we establish regret bounds of order $T^{(\beta+1)/(\beta+2)}$, where $\beta > 0$ is a hyperparameter that balances reward maximization against mitigating distribution shift. To our knowledge, this is the first online RLHF algorithm with regret scaling polynomially in all model parameters.

1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable capabilities across a wide range of natural language tasks, yet aligning their behavior with human preferences remains a central challenge. A widely adopted solution is reinforcement learning with human feedback (RLHF), which fine-tunes a pretrained LLM using human preference data (Christiano et al., 2017; Ziegler et al., 2019; Bai et al., 2022). The standard RLHF pipeline involves three stages: (i) supervised fine-tuning (SFT) on human-written demonstrations to produce a baseline model; (ii) training a reward model from human preference comparisons (Bradley & Terry, 1952); and (iii) optimizing the LLM with reinforcement learning against the learned reward. This framework has been instrumental in the success of instruction-following LLMs such as InstructGPT (Ouyang et al., 2022) and ChatGPT (OpenAI, 2023), enabling models to produce responses that are more helpful, safe, and aligned with human expectations.

Despite this progress, most existing RLHF implementations are offline (Zhao et al., 2023; Rafailov et al., 2024; Azar et al., 2024): the preference data is collected once from static policies, and the reward model is trained on this fixed dataset (Ivison et al., 2023; Zhu et al., 2024; Shi et al., 2025). While effective, offline RLHF has inherent limitations—It cannot adaptively explore the enormous space of natural language, leading to inefficient use of expensive human feedback. In contrast, online RLHF offers a more powerful alternative: the policy iteratively collects new preference data, updates the reward model, and improves itself based on these updates (Guo et al., 2024; Xiong et al., 2023; Chen et al., 2024; Rosset et al., 2024; Dong et al., 2024; Feng et al., 2025). This interactive loop has the potential to greatly improve both alignment quality and sample efficiency. However, realizing this potential requires principled approaches to exploration, i.e., deciding which comparisons to query in order to most effectively reduce uncertainty in reward estimation.

054 A natural candidate for encouraging and guiding exploration is the principle of optimism (Lai &
 055 Robbins, 1985; Lattimore & Szepesvári, 2020), which acts as if the environment is more optimistic
 056 than currently estimated, within the limits of statistical uncertainty based on all data that has been
 057 observed so far. It is usually implemented by adding an uncertainty-based bonus to reward or
 058 value estimates, thereby prioritizing actions whose values are uncertain but potentially high. This
 059 has yielded provably efficient algorithms in standard RL (see e.g., Jin et al. (2018); Zanette &
 060 Brusiloff (2019); Russo & Van Roy (2013); Azar et al. (2017)). However, extending this principle to
 061 RLHF introduces new difficulties, where feedback comes not as a single reward but as a difference
 062 between rewards of two actions. The key challenge is to determine the action pairs with the large
 063 uncertainties most relevant to policy improvement. A few recent works achieved important progress
 064 towards designing sample-efficient online RLHF algorithms based on the optimism principle (Cen
 065 et al., 2025; Zhang et al., 2025; Xie et al., 2025). However the existing theoretical guarantees still
 066 exhibit exponential dependency on certain model parameters, which potentially leads to inefficient
 067 exploration.

068 With this context, this paper makes contribution towards designing efficient online exploration
 069 schemes for RLHF with provable guarantees. By analyzing the existing algorithms in the seminal
 070 works (Cen et al., 2025; Xie et al., 2025; Zhang et al., 2025), we discuss their inadequacy
 071 in exploring the action pairs with the large uncertainties most relevant to policy improvement,
 072 and construct lower bounds to show that the exponential dependency on certain parameters is
 073 unavoidable in their regret. Based on these insights, we propose a new exploration scheme for RLHF
 074 that adopts a different sampling protocol, and establish a regret bound that depends polynomially on
 075 all model parameters.

2 MODEL SET-UP

076 **Preliminaries.** In RLHF, the prompt space \mathcal{X} refers to the collection of all possible inputs or
 077 queries that a user might provide to the model. The answer (or action) space \mathcal{A} is the set of all
 078 possible outputs the model can generate in reply to a given prompt. A language model is a policy
 079 $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$ that defines a probability distribution $\pi(\cdot | x)$ over \mathcal{A} conditioned on a prompt
 080 $x \in \mathcal{X}$, specifying how likely the model is to produce each potential response. The pipeline of
 081 RLHF starts with supervised fine-tuning (SFT), where a reference policy $\pi_{\text{ref}} : \mathcal{X} \rightarrow \Delta(\mathcal{A})$ is
 082 obtained by fine-tuning a pre-trained LLM on a dataset of prompts paired with high-quality answers
 083 written by humans. SVT provides an initialization that stabilizes and improves the effectiveness of
 084 the subsequent training stages that aligns the LLM with human preferences.

085 **Reward modeling.** To translate human preferences into a trainable objective, one need to model
 086 how an oracle (e.g., a human annotator) rank two answers a_1 and a_2 given prompt x . Following a
 087 line of prior works (e.g., Cen et al. (2025); Xie et al. (2025); Zhang et al. (2025)), we assume that
 088 preferences follow the Bradley-Terry model (Bradley & Terry, 1952)

$$089 \mathbb{P}(a_1 \succ a_2 | x) = \frac{\exp(r^*(x, a_1))}{\exp(r^*(x, a_1)) + \exp(r^*(x, a_2))} = \sigma(r^*(x, a_1) - r^*(x, a_2)). \quad (2.1)$$

090 Here $r^* : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ is an underlying reward function of an answer given a prompt, $a_1 \succ a_2$
 091 means the answer a_1 is preferred compared to a_2 , and $\sigma(x) = (1 + e^{-x})^{-1}$ is the sigmoid function.
 092 We also define a policy π_{HF} to characterize human preference:

$$093 \pi_{\text{HF}}(a | x) = \frac{\exp(r^*(x, a))}{\sum_{a' \in \mathcal{A}} \exp(r^*(x, a'))}.$$

094 The reward function is unknown and can be learned from e.g., an offline dataset $\mathcal{D} = \{(x^i, a_+^i, a_-^i)\}$
 095 comprised of independent preference data samples using maximum likelihood estimation (MLE):

$$096 \arg \max_r \ell(r, \mathcal{D}) \quad \text{where} \quad \ell(r, \mathcal{D}) := \sum_{\mathcal{D}} \log \sigma(r(x^i, a_+^i) - r(x^i, a_-^i)), \quad (2.2)$$

097 where a preference data sample denoted by (x, a_+, a_-) means that $a_+ \succ a_-$ given prompt x .

108 **RL fine-tuning.** Given a reward model r , we seek to fine-tune the policy π to balance reward
 109 maximization with maintaining similarity to the original model π_{ref} from the SFT stage. Towards
 110 this, we define the KL-regularized reward objective
 111

$$112 J(\pi, r; \pi_{\text{cal}}) := \mathbb{E}_{x \sim \rho} [\mathbb{E}_{a \sim \pi(\cdot | x)} [r(x, a)] - \mathbb{E}_{a \sim \pi_{\text{cal}}(\cdot | x)} [r(x, a)] - \beta \text{KL}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]. \quad (2.3)$$

113 Here ρ is the prompt distribution, and $\beta > 0$ is the regularization parameter reflecting the strength
 114 of the KL regularization. In practice, β is typically chosen to be small; for instance, in InstructGPT
 115 (Ouyang et al., 2022) the optimal value is reported to be around 0.01 and 0.02. This objective
 116 function includes a calibration policy π_{cal} to eliminate the shift ambiguity of the reward function, as
 117 two reward functions $r(x, a)$ and $r(x, a) + c(x)$ lead to the same preference model (2.1). Given any
 118 reward function r , the optimal policy $\pi_r := \arg \max_{\pi} J(\pi, r; \pi_{\text{cal}})$ admits a closed-form expression
 119 (Rafailov et al., 2024)

$$120 \pi_r(a | x) = \frac{\pi_{\text{ref}}(a | x) \exp(r(x, a) / \beta)}{Z_r(x)} \quad (2.4)$$

121 where $Z_r(x) = \sum_a \pi_{\text{ref}}(a | x) \exp(r(x, a) / \beta)$ is the normalizing factor. Notice that the selection
 122 of π_{cal} does not affect the optimal policy π_r given the reward function r . Our target is the optimal
 123 policy π^* that maximizes the objective (2.3) under the true reward function $r = r^*$, namely
 124

$$125 \pi^* := \arg \max_{\pi} J(\pi, r^*; \pi_{\text{cal}}). \quad (2.5)$$

126 **Offline RLHF.** The above framework leads to offline RLHF methods that relies on the preference
 127 dataset \mathcal{D} for training. Initial approaches (Christiano et al., 2017; Ouyang et al., 2022) first
 128 estimate a reward function \hat{r} based on the preference dataset \mathcal{D} using MLE, then optimize the KL-
 129 regularized objective (2.3) with respect to \hat{r} . Another approach introduced by Rafailov et al. (2024)
 130 condensed these two steps into one single step, known as direct preference optimization (DPO),
 131 which optimizes
 132

$$133 \max_{\pi} \sum_{\mathcal{D}} \log \sigma \left(\beta \left(\log \frac{\pi(y_+^i | x)}{\pi_{\text{ref}}(y_+^i | x)} - \log \frac{\pi(y_-^i | x)}{\pi_{\text{ref}}(y_-^i | x)} \right) \right).$$

134 The above objective avoids explicitly estimating the reward function, which can be obtained by
 135 expressing the reward function r in the MLE formulation (2.2) with the associated optimal policy
 136 π_r using the closed-form expression (2.4). However, as discussed in e.g., Xie et al. (2025); Zhang
 137 et al. (2025), the efficiency of offline RLHF is limited by the coverage of the offline dataset \mathcal{D} , and
 138 online exploration with active data collection is necessary to achieve sample efficiency.

139 **Online RLHF.** We consider reward learning and policy learning iteratively, where in the t -th
 140 iteration we use the current policy $\pi^{(t)}$, obtained from previous iterations, to sample new data and
 141 subsequently update both the reward estimate and the policy. This setup enables online exploration
 142 in RLHF by refining the reward model and policy in tandem as new preference data is collected. We
 143 aim to minimize the regret

$$144 \mathcal{R}(T) := \sum_{t=1}^T [J(\pi^*; r^*, \pi_{\text{cal}}) - J(\pi^{(t)}; r^*, \pi_{\text{cal}})]. \quad (2.6)$$

145 It is worth mentioning that the choice of π_{cal} does not affect the regret. We define the following
 146 function J^* that measures the optimal objective value for a given reward r :

$$147 J^*(r; \pi_{\text{cal}}) := \max_{\pi} J(\pi, r; \pi_{\text{cal}}) = J(\pi_r, r; \pi_{\text{cal}}). \quad (2.7)$$

148 This function plays an important role in the exploration algorithms.

149 3 RLHF WITH ONLINE EXPLORATION

150 Three recent algorithms for online RLHF are most closely related to this work: VPO (Cen et al.,
 151 2025), XPO (Xie et al., 2025), and SELM (Zhang et al., 2025). In this section, we first analyze and
 152 discuss these approaches, and then introduce our proposed exploration scheme.

162 3.1 INADEQUACY OF EXISTING APPROACHES
163

164 We begin by reviewing the procedure and intuition behind VPO (Cen et al., 2025). Fix a calibration
165 policy π_{cal} and an initial policy $\pi^{(1)}$. For $t = 1, 2, \dots, T$, the t -th iteration of VPO consists of the
166 following steps:

167 1. Sample a prompt $x^t \sim \rho$ and two answers $a_1^t, a_2^t \sim \pi^{(t)}(\cdot | x^t)$. Query the preference oracle
168 to obtain pairwise comparison $a_+^t \succ a_-^t$. Update the preference dataset $\mathcal{D}^{(t)} = \mathcal{D}^{(t-1)} \cup$
169 $\{(x^t, a_+^t, a_-^t)\}$.

170 2. Update the reward model $r^{(t+1)}$ and the policy $\pi^{(t+1)}$ using the updated preference dataset $\mathcal{D}^{(t)}$:

$$r^{(t+1)} = \arg \max_{r: \mathcal{X} \times \mathcal{A} \rightarrow [0, r_{\max}]} \ell(r, \mathcal{D}^{(t)}) + \alpha J^*(r; \pi_{\text{cal}}), \quad (3.1a)$$

$$\pi^{(t+1)} = \arg \max_{\pi} J(\pi, r^{(t+1)}; \pi_{\text{cal}}), \quad (3.1b)$$

171 where $\alpha > 0$ is a regularization parameter, and step (3.1b) admits closed-form solution (2.4).

172 To illustrate the rationale behind VPO, consider the bandit case with no prompt. Step (3.1a) applies
173 the optimism principle, encouraging exploration based on the uncertainty in estimating the reward
174 difference between each action a and the calibration policy π_{cal} . Formally, it can be viewed as the
175 Lagrangian form of the constrained optimization problem

$$\max_{r, \pi} \mathbb{E}_{a \sim \pi}[r(a)] - \mathbb{E}_{a \sim \pi_{\text{cal}}}[r(a)] - \beta \text{KL}(\pi \| \pi_{\text{ref}}) \quad \text{s.t.} \quad \ell(r, \mathcal{D}^{(t)}) \geq \max_r \ell(r, \mathcal{D}^{(t)}) - B$$

176 for some $B > 0$. After the change of variable $r'(a) = r(a) - \mathbb{E}_{a \sim \pi_{\text{cal}}}[r(a)]$, this becomes

$$\max_{r', \pi} \mathbb{E}_{a \sim \pi}[r'(a)] - \beta \text{KL}(\pi \| \pi_{\text{ref}}) \quad \text{s.t.} \quad \ell(r', \mathcal{D}^{(t)}) \geq \max_{r'} \ell(r', \mathcal{D}^{(t)}) - B, \quad \mathbb{E}_{a \sim \pi_{\text{cal}}}[r'(a)] = 0.$$

177 Here, the constraint set can be interpreted as a confidence region reflecting the uncertainty in
178 estimating each $r'(a)$ from $\mathcal{D}^{(t)}$. Consequently, the updated policy $\pi^{(t+1)}$ depends both on the
179 true reward gap $r(a) - \mathbb{E}_{a \sim \pi_{\text{cal}}}[r(a)]$ and on the uncertainty in estimating this gap for each action
180 $a \in \mathcal{A}$.

181 For intuition, suppose $\pi_{\text{cal}} = \mathbb{1}_{a_0}$ for some $a_0 \in \mathcal{A}$, and assume that the true reward gaps are small.
182 In this case, $\pi^{(t+1)}$ favors actions with higher estimation uncertainty relative to a_0 , i.e., those a where
183 the estimate of $r(a) - r(a_0)$ is most uncertain. However, comparing two actions $a_1, a_2 \sim \pi^{(t+1)}$
184 reduces the uncertainty between them, rather than the (potentially larger) uncertainty relative to a_0 .
185 This misalignment can lead to inefficient exploration, as illustrated in the following example.

186 **Example 1.** Consider the bandit setting with three actions $\mathcal{A} = \{a_0, a_1, a_2\}$, where the true rewards
187 are $r^*(a_0) = 1$ and $r^*(a_1) = r^*(a_2) = 0$. Let the reference policy π_{ref} be uniform over \mathcal{A} , and the
188 calibration policy be $\pi_{\text{cal}}(a_1) = \pi_{\text{cal}}(a_2) = p$ and $\pi_{\text{cal}}(a_0) = 1 - 2p$ for some $0 \leq p < 1/4$.

189 The following proposition shows that VPO may fail to explore efficiently in this setting. The proof
190 can be found in Appendix A.

191 **Proposition 1.** Consider the setup in Example 1. Let the initial policy $\pi^{(1)}$ of VPO be the uniform
192 distribution over \mathcal{A} . Assume that $r_{\max}/\beta \geq 3$. For any $\alpha > 0$, with probability at least $4/(9e)$, we
193 have

$$J(\pi^*, r^*; \pi_{\text{cal}}) - J(\pi^{(t)}, r^*; \pi_{\text{cal}}) \geq \frac{1}{2}$$

194 holds for any $1 < t \leq \exp(r_{\max}/\beta)/2$.

195 Let's discuss the idea behind Proposition 1 with $\pi_{\text{cal}} = \mathbb{1}_{a_0}$. If the calibration action a_0 is
196 not visited during the first t iterations, then $\pi^{(t+1)}$ will continue to favor a_1 and a_2 , since both
197 gaps $r(a_1) - r(a_0)$ and $r(a_2) - r(a_0)$ remain highly uncertain. In particular, we establish that
198 $\pi^{(t+1)}(a_0) \leq \exp(-r_{\max}/\beta)$, which is exponentially small, implying that a_0 is unlikely to be
199 sampled in iteration $t + 1$. As a result, with constant probability, a_0 will not be sampled within
200 the first $O(\exp(r_{\max}/\beta))$ iterations, and the resulting highly suboptimal policy incurs linear regret
201 over an exponentially long horizon. This example highlights an algorithmic drawback: although
202 VPO acknowledges uncertainty in the reward gaps between a_1 and a_0 (and between a_2 and a_0), it
203 continues to encourage sampling a_1 and a_2 , leading primarily to comparisons between them that fail
204 to reduce their uncertainty relative to a_0 .

216
217**Algorithm 1:** Uncertainty-based RLHF exploration.218
2191 **Input:** initial policies $\pi^{(0)}, \pi^{(1)}$, regularizaton parameters $\{\alpha_t\}_{t \geq 1}$.2 **for** $t = 1$ **to** T **do**

220

3 Sample a prompt $x^t \sim \rho$ and two answers $a_1^t \sim \pi^{(t-1)}(\cdot | x^t), a_2^t \sim \pi^{(t)}(\cdot | x^t)$.

221

4 Query the preference oracle to obtain pairwise comparison $a_+^t \succ a_-^t$ and update the preference dataset $\mathcal{D}^{(t)} = \mathcal{D}^{(t-1)} \cup \{(x^t, a_+^t, a_-^t)\}$.

222

5 Update the reward model $r^{(t+1)}$ and the policy $\pi^{(t+1)}$ using $\mathcal{D}^{(t)}$:

223

$$r^{(t+1)} = \arg \max_{r: \mathcal{X} \times \mathcal{A} \rightarrow [0, r_{\max}]} \ell(r, \mathcal{D}^{(t)}) + \alpha_t J^*(r; \pi^{(t)}), \quad (3.2a)$$

224

$$\pi^{(t+1)} = \arg \max_{\pi} J(\pi, r^{(t+1)}; \pi^{(t)}). \quad (3.2b)$$

225

where the policy update (3.2b) admits closed-form solution (2.4).

226

6 **Output:** $\{\pi^{(t)} : 1 \leq t \leq T\}$

227

228

3.2 OUR APPROACH: EXPLORATION BASED ON UNCERTAINTY

229

A natural modification to address the issue above is to change the sampling scheme so that $a_1^t \sim \pi^{(t)}$ and $a_2^t \sim \pi_{\text{cal}}$. The intuition is that $\pi^{(t)}$ encourages to explore actions with higher estimation uncertainty relative to the actions favored by the calibration policy π_{cal} . To effectively reduce this uncertainty, it is sensible to compare one action drawn from $\pi^{(t)}$ with another drawn from π_{cal} . Indeed, the XPO and SELM algorithms (Xie et al., 2025; Zhang et al., 2025) can be viewed as taking $\pi_{\text{cal}} = \pi_{\text{ref}}$.

230

However, if the fixed calibration policy π_{cal} is highly suboptimal for reward maximization (for example, if it concentrates on a few low-reward actions), then the comparison will almost always favor $a_1^t \sim \pi^{(t)}$ against $a_2^t \sim \pi_{\text{cal}}$, yielding little useful information. This issue is illustrated in the following example.

231

232

233

234

Example 2. Consider the bandit setting with three actions $\mathcal{A} = \{a_0, a_1, a_2\}$, where the true rewards are $r^*(a_0) = 0, r^*(a_1) = r_{\max}$ and $r^*(a_2) = r_{\max} - 2$. Let the reference policy be $\pi_{\text{ref}}(a_0) = 1 - 2/\kappa, \pi_{\text{ref}}(a_1) = \pi_{\text{ref}}(a_2) = 1/\kappa$ for any $\kappa \geq 4$.

235

The following result shows that, when κ is large (as we will see in Assumption 1, this corresponds to the case where the reference policy deviates from human preference), this modified sampling schemes can lead to inefficient exploration in this setting. The proof is deferred to Appendix B.

236

Proposition 2. Consider the setup in Example 2. Assume that $\beta \leq 1$ and $\kappa \leq \exp(r_{\max}/\beta)$. For any initial policy $\pi^{(1)}$ and any $\alpha > 0$, with probability at least $1/64$, the modified exploration scheme which samples $a_1^t \sim \pi^{(t)}$ and $a_2^t \sim \pi_{\text{ref}}$ satisfies

237

$$J(\pi^*, r^*; \pi_{\text{ref}}) - J(\pi^{(t)}, r^*; \pi_{\text{ref}}) \geq 0.01$$

238

for any $1 < t \leq \min\{\kappa, \exp(r_{\max})/2\}$.

239

This lower bound suggests that relying on a fixed calibration policy can lead to inefficient exploration over an exponentially long horizon. We will come back to this example in Section 4 after presenting our algorithm and theoretical guarantees. This observation motivates us to update the calibration policy in each iteration adaptively.

240

241

242

Uncertainty-based exploration. We propose an exploration scheme where the calibration policy evolves with the iterations. In the t -th iteration, instead of a fixed π_{cal} , we use $\pi^{(t)}$ as the calibration policy when optimizing $r^{(t+1)}$ and $\pi^{(t+1)}$:

243

244

245

$$r^{(t+1)} = \arg \max_{r: \mathcal{X} \times \mathcal{A} \rightarrow [0, r_{\max}]} \ell(r, \mathcal{D}^{(t)}) + \alpha_t J^*(r; \pi^{(t)}),$$

246

247

248

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

249

270 The key advantage is that $\pi^{(t)}$ improves over time, guiding exploration away from uninformative
 271 comparisons. Since $\pi^{(t)}$ emphasizes actions with higher uncertainty relative to $\pi^{(t-1)}$, it is natural
 272 to compare $a_1^t \sim \pi^{(t-1)}$ and $a_2^t \sim \pi^{(t)}$. This yields preference data that more directly reduces
 273 uncertainty, leading to more efficient exploration. Our full exploration scheme is summarized in
 274 Algorithm 1.

276 4 THEORETICAL RESULTS

278 We establish theoretical guarantees for Algorithm 1 under the multi-armed bandit setting (i.e., $\mathcal{X} =$
 279 \emptyset) with $A = |\mathcal{A}|$. We begin with a general regret bound, whose proof is deferred to Section 5.

280 **Theorem 1.** *Let $\alpha_t > A \log T$ be non-decreasing in t . There exists a universal constant $C > 0$
 281 such that, with probability at least $1 - O(T^{-10})$, the cumulative regret of running Algorithm 1 for
 282 T iterations satisfies*

$$284 \mathcal{R}(T) \leq Cr_{\max}A^2\sqrt{T \log T} + C \sum_{t=1}^T \frac{Ar_{\max} \log T}{\alpha_t} + CA^2\alpha_T r_{\max}^2 \quad (4.1) \\ 285 \\ 286 \\ 287 + C(r_{\max} + \log T) \sum_{r^*(a_+) \geq r^*(a_-)} \min \left\{ \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T r_{\max}, \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\}. \\ 288 \\ 289$$

290 We now discuss the implications of Theorem 1. When $\beta = 0$, which corresponds to the case where
 291 only reward maximization matters, the regret bound (4.1) simplifies to

$$293 \mathcal{R}(T) = \tilde{O}((A^{3/2}r_{\max}^{3/2} + A^2r_{\max})\sqrt{T}) \quad \text{when} \quad \alpha_t \asymp A \log T + \sqrt{\frac{t}{Ar_{\max}}}. \\ 294$$

295 When $\beta > 0$, the performance of the exploration algorithm becomes more intricate due to the
 296 trade-off between reward maximization and similarity to the reference policy. To interpret the
 297 general regret bound in this regime, we introduce the following assumption to capture the interaction
 298 between human preference π_{HF} and the reference policy π_{ref} .

299 **Assumption 1.** *There exists $\kappa, \tau \geq 1$ such that, for any action pair (a_+, a_-) ,*

$$300 \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \geq \tau \implies \frac{\pi_{\text{ref}}(a_+)}{\pi_{\text{ref}}(a_-)} \geq \kappa^{-1}. \\ 301 \\ 302$$

303 Intuitively, Assumption 1 requires that whenever a_+ is substantially more preferred than a_- under
 304 human preference, the reference policy does not assign disproportionately higher weight to a_- than
 305 to a_+ . This is reasonable, since π_{ref} is obtained from the SFT step, where a pretrained LLM is
 306 fine-tuned on human demonstrations already broadly aligned with preference. The quantities κ and
 307 τ capture the degree of alignment between π_{ref} and π_{HF} , and their size reflects the influence of
 308 the reference policy on RLHF. We note that the illustrative Example 1 satisfies Assumption 1 with
 309 $\kappa, \tau = O(1)$, and the parameter κ in Example 2 is consistent with the κ here. Under this assumption,
 310 we obtain the following simplified regret bound, whose proof is deferred to Appendix D.

311 **Proposition 3.** *Suppose that Assumption 1 holds. Let*

$$312 \alpha_t = A \log T + t^{\frac{1}{\beta+2}} \left(\frac{r_{\max}}{\kappa} \right)^{\frac{\beta}{\beta+2}} \left(\frac{\log T}{A(r_{\max} + \log T)} \right)^{\frac{\beta+1}{\beta+2}}. \\ 313$$

314 Then with probability at least $1 - O(T^{-10})$, we have

$$315 \mathcal{R}(T) \lesssim (\tau + \kappa^\beta T^{\frac{\beta+1}{\beta+2}}) \text{poly}(A, r_{\max}, \log T),$$

317 where the degree of the polynomial factor does not depend on β .

318 **Remark 1.** When κ is large, namely the reference policy deviates significantly from the human
 319 preference, it is natural to choose a small KL regularization parameter β to reduce the influence of
 320 the reference policy. In this regime, Algorithm 1 remains robust, since the regret bound scales only
 321 with κ^β . By contrast, the lower bound in Proposition 2 suggests that the sampling protocols in prior
 322 works (Xie et al., 2025; Zhang et al., 2025) would incur regret at least linear in κ . This demonstrates
 323 that our strategy accommodates scenarios with small β , where the reference policy is poorly aligned
 with human preference.

324 **Remark 2.** In Appendix E, we present an alternative assumption linking human preference and the
 325 reference policy, together with the corresponding regret guarantee.
 326

327 Proposition 3 establishes a regret bound of order $O(T^{\frac{\beta+1}{\beta+2}})$, with only polynomial dependence on
 328 the other parameters. This stands in sharp contrast to prior works (Cen et al., 2025; Xie et al., 2025;
 329 Zhang et al., 2025), which achieved the more standard $O(\sqrt{T})$ regret but at the cost of exponential
 330 dependence on terms such as r_{\max}/β . We conjecture that, for RLHF, eliminating exponential
 331 dependence inevitably requires a slower rate in T , with the exponent governed by β . This trade-off
 332 is intuitive: online exploration primarily serves to learn human preference, and as the regularization
 333 parameter β increases, greater emphasis is placed on preserving similarity to the reference measure.
 334 This constraint naturally slows convergence.

335 5 PROOF OF THEOREM 1

336 5.1 STEP 1: REGRET DECOMPOSITION

337 In view of the optimality of $r^{(t)}$ (cf. equation (3.2a)), we have

$$338 \quad \ell(r^{(t)}, \mathcal{D}^{(t-1)}) + \alpha_t J^*(r^{(t)}; \pi^{(t-1)}) \geq \ell(r^*, \mathcal{D}^{(t-1)}) + \alpha_t J^*(r^*; \pi^{(t-1)}).$$

339 Rearrange terms to get

$$340 \quad \begin{aligned} \frac{1}{\alpha_t} [\ell(r^{(t)}, \mathcal{D}^{(t-1)}) - \ell(r^*, \mathcal{D}^{(t-1)})] &\geq J^*(r^*; \pi^{(t-1)}) - J^*(r^{(t)}; \pi^{(t-1)}) \\ &\stackrel{(i)}{=} \max_{\pi} J(\pi, r^*; \pi^{(t-1)}) - \max_{\pi} J(\pi, r^{(t)}; \pi^{(t-1)}) \\ &\stackrel{(ii)}{\geq} J(\pi^*, r^*; \pi^{(t-1)}) - J(\pi^{(t)}, r^{(t)}; \pi^{(t-1)}). \end{aligned} \quad (5.1)$$

341 Here step (i) follows from the definition of J^* (cf. equation (2.7)), while step (ii) follows from the
 342 optimality of $\pi^{(t)}$ (cf. equation (3.2b)). This allows us to reach the following decomposition:

$$343 \quad \begin{aligned} \text{Regret}_t &:= J(\pi^*, r^*; \pi^{(t-1)}) - J(\pi^{(t)}, r^*; \pi^{(t-1)}) \\ &\leq \underbrace{\alpha_t^{-1} [\ell(r^{(t)}, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)})]}_{=: \theta_t} + \underbrace{J(\pi^{(t)}, r^{(t)}; \pi^{(t-1)}) - J(\pi^{(t)}, r^*; \pi^{(t-1)})}_{=: \gamma_t}. \end{aligned} \quad (5.2)$$

344 In view of the definition of J (cf. equation (2.3)), we can further decompose

$$345 \quad \begin{aligned} \gamma_t &= \mathbb{E}_{a \sim \pi^{(t)}} [r^{(t)}(a)] - \mathbb{E}_{a \sim \pi^{(t-1)}} [r^{(t)}(a)] - \mathbb{E}_{a \sim \pi^{(t)}} [r^*(a)] + \mathbb{E}_{a \sim \pi^{(t-1)}} [r^*(a)] \\ &= r^{(t)}(a_2^t) - r^{(t)}(a_1^t) - r^*(a_2^t) + r^*(a_1^t) + \xi_t \end{aligned}$$

346 where ξ_t is the martingale difference sequence

$$347 \quad \begin{aligned} \xi_t &= \mathbb{E}_{a \sim \pi^{(t)}} [r^{(t)}(a)] - r^{(t)}(a_2^t) - \mathbb{E}_{a \sim \pi^{(t-1)}} [r^{(t)}(a)] + r^{(t)}(a_1^t) \\ &\quad - \mathbb{E}_{a \sim \pi^{(t)}} [r^*(a)] + r^*(a_2^t) + \mathbb{E}_{a \sim \pi^{(t-1)}} [r^*(a)] - r^*(a_1^t). \end{aligned}$$

348 Therefore we have

$$349 \quad \text{Regret} = \sum_{t=1}^T \text{Regret}_t \leq \underbrace{\sum_{t=1}^T \theta_t}_{=: \theta} + \underbrace{\sum_{t=1}^T \xi_t}_{=: \xi} + \underbrace{\sum_{t=1}^T |r^{(t)}(a_2^t) - r^{(t)}(a_1^t) - r^*(a_2^t) + r^*(a_1^t)|}_{=: \zeta}. \quad (5.3)$$

350 It is straightforward to bound the second term ξ . Notice that $|\xi_t| \leq 8r_{\max}$ holds deterministically
 351 for any $1 \leq t \leq T$. By the Azuma-Hoeffding inequality, with probability exceeding $1 - O(T^{-10})$
 352 we have

$$353 \quad \xi = \sum_{t=1}^T \xi_t \leq C_1 r_{\max} \sqrt{T \log T} \quad (5.4)$$

354 for some universal constant $C_1 > 0$. In what follows, we bound the other two terms θ and ζ .

378 5.2 STEP 2: BOUNDING LIKELIHOOD RATIOS
379380 To bound θ , we need to analyze the regularized MLE. Notice that
381

382
$$\theta_t = \frac{\ell(r^{(t)}, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)})}{\alpha_t} = \alpha_t^{-1} \sum_{i=1}^t \log \frac{\sigma(r^{(t)}(x^i, a_+^i) - r^{(t)}(x^i, a_-^i))}{\sigma(r^*(x^i, a_+^i) - r^*(x^i, a_-^i))}.$$

383

384 The following lemma is crucial for the subsequent analysis. The proof can be found in Appendix C.1.
385386 **Lemma 1.** For any given reward function $r : \mathcal{A} \rightarrow [\pm r_{\max}]$ and any $1 \leq t \leq T$, define
387

388
$$\Delta_t(r) := \sum_{i=1}^t \log \frac{\sigma(r^*(a_+^i) - r^*(a_-^i))}{\sigma(r(a_+^i) - r(a_-^i))} - \sum_{i=1}^t \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))).$$

389

390 There exists some universal constant $C_2 > 1$ such that for any fixed r , with probability at least $1 - \delta$,
391

392
$$|\Delta_t(r)| \leq C_2 \sqrt{\sum_{i=1}^t r_{\max} \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) \log \frac{\log T}{\delta} + C_2 r_{\max} \log \frac{\log t}{\delta}}.$$

393

394 Equipped with the concentration bounds in Lemma 1, we can use the standard covering argument to
395 derive an uniform upper bound, whose proof is deferred to Appendix C.2.
396397 **Lemma 2.** There exists some universal constant $C_3 > 0$ such that with probability exceeding $1 - O(T^{-9})$,
398

399
$$\ell(r, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)}) \leq -\frac{1}{2} \sum_{i=1}^t \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) + C_3 A r_{\max} \log T$$

400

401 holds for any $r : \mathcal{A} \rightarrow [\pm r_{\max}]$ and $1 \leq t \leq T$.
402403 As an immediate consequence of Lemma 2, with probability exceeding $1 - O(T^{-9})$,
404

405
$$\ell(r^{(t)}, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)}) \leq C_3 A r_{\max} \log T$$

406

407 holds for any $1 \leq t \leq T$. Therefore
408

409
$$\theta = \sum_{t=1}^T \theta_t = \sum_{t=1}^T \frac{\ell(r^{(t)}, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)})}{\alpha_t} \leq \sum_{t=1}^T \frac{C_3 A r_{\max} \log T}{\alpha_t}. \quad (5.5)$$

410

412 5.3 STEP 3: BOUNDING REWARD ERRORS
413414 We first notice that
415

416
$$\begin{aligned} \alpha_t^{-1} [\ell(r^{(t)}, \mathcal{D}^{(t-1)}) - \ell(r^*, \mathcal{D}^{(t-1)})] &\stackrel{(i)}{\geq} \max_{\pi} J(\pi, r^*; \pi^{(t-1)}) - \max_{\pi} J(\pi, r^{(t)}; \pi^{(t-1)}) \\ &\stackrel{(ii)}{\geq} J(\pi^{(t)}, r^*; \pi^{(t-1)}) - J(\pi^{(t)}, r^{(t)}; \pi^{(t-1)}) \\ &\stackrel{(iii)}{=} \mathbb{E}_{a \sim \pi^{(t)}} [r^{(t)}(a) - r^*(a)] - \mathbb{E}_{a \sim \pi^{(t-1)}} [r^{(t)}(a) - r^*(a)] \\ &\geq -4r_{\max}. \end{aligned} \quad (5.6)$$

417

418 Here step (i) is an intermediate step of (5.1); step (ii) follows from the optimality of $\pi^{(t)}$ (cf. (3.2b));
419 step (iii) follows from the definition of J (cf. (2.3)). This combined with Lemma 2 implies that
420

421
$$\begin{aligned} \sum_{i=1}^t \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r^{(t)}(a_1^i) - r^{(t)}(a_2^i))) \\ \leq -2[\ell(r^{(t)}, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)})] + 2C_3 A r_{\max} \log T \leq C_4 \alpha_t r_{\max}, \end{aligned} \quad (5.7)$$

422

423 as long as $\alpha_t \geq A \log T$ and $C_4 \geq 8 + 2C_3$. This implies that for any $t \in [T]$ and any action pair
424 (a_+, a_-) ,
425

426
$$\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t)}(a_+) - r^{(t)}(a_-))) \leq \frac{C_4 \alpha_t r_{\max}}{N_t(a_+, a_-)}, \quad (5.8)$$

427

432 where $N_t(a_+, a_-)$ is the number of comparison for (a_+, a_-) up to time t . This motivates us to
 433 decompose ζ according to whether $N_t(a_+, a_-) \gg \alpha_t r_{\max}$: let $\tau := 100C_4\alpha_T r_{\max}$ and denote by
 434 $t_n(a_+, a_-)$ the time of the n -th comparison for (a_+, a_-) , we have

$$436 \quad \zeta \leq 2\tau A^2 r_{\max} + \sum_{r^*(a_+) \geq r^*(a_-)} \underbrace{\sum_{n=\tau}^{N_T(a_+, a_-)} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)|}_{=: \zeta(a_+, a_-)},$$

440 where we denote by $t_n(a_+, a_-)$ the time of the n -th comparison for (a_+, a_-) , and the first
 441 summation is taken over all action pairs (a_+, a_-) satisfying $r^*(a_+) \geq r^*(a_-)$. To bound each
 442 $\zeta(a_+, a_-)$, we need the following technical lemma. The proof can be found in Appendix C.3.

443 **Lemma 3.** *Consider any action pair (a_+, a_-) and time t_0 such that $N_{t_0}(a_+, a_-) \geq \tau$. There
 444 exists universal constant $C_5 > 0$ such that, for any $t_0 \leq t_1 < t_2 \leq T$, with probability exceeding
 445 $1 - O(T^{-10})$ we have*

$$447 \quad N_{t_2}(a_+, a_-) - N_{t_1}(a_+, a_-) \leq C_5 \sum_{t=t_1+1}^{t_2} \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \left[\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \| \sigma(r^{(t)}(a_+) - r^{(t)}(a_-)))^{\frac{1}{\beta}} \right. \\ 448 \quad \left. + \sigma(r^*(a_-) - r^*(a_+))^{\frac{1}{\beta}} \right] + C_5 \sqrt{T \log T}.$$

451 Equipped with Lemma 3, we can bound each $\zeta(a_+, a_-)$ using both density ratios regarding human
 452 feedback $\pi_{\text{HF}}(a_+)/\pi_{\text{HF}}(a_-)$, and regarding the reference policy $\pi_{\text{ref}}(a_-)/\pi_{\text{ref}}(a_+)$. The proof is
 453 deferred to Appendix C.4.

454 **Lemma 4.** *There exists universal constant $C_6 > 0$ such that, for any action pair (a_+, a_-) , with
 455 probability exceeding $1 - O(T^{-9})$ we have*

$$457 \quad \zeta(a_+, a_-) \leq C_6(r_{\max} + \log T) \min \left\{ \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T r_{\max}, \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\} \\ 458 \quad + C_6 \left(\frac{AN_T(a_+, a_-) \log T}{\alpha_T} + \sqrt{T \log T} \right) r_{\max}.$$

462 This immediately implies that

$$464 \quad \zeta \leq 2\tau A^2 r_{\max} + C_6 \left(\frac{AT \log T}{\alpha_T} + A^2 \sqrt{T \log T} \right) r_{\max} \quad (5.9) \\ 465 \quad + C_6(r_{\max} + \log T) \sum_{r^*(a_+) \geq r^*(a_-)} \min \left\{ \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T r_{\max}, \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\}.$$

470 Putting the regret decomposition (5.3) and the bounds (5.4), (5.5) and (5.9) collectively yields the
 471 desired regret bound (4.1).

472 6 DISCUSSION

475 In this paper, we investigated the problem of efficient exploration in online RLHF. By a careful
 476 analysis of the existing optimism-based exploration strategies, we identified a conceptual drawback
 477 in their sampling protocol, and we proved lower bounds to show that they can lead to inefficient
 478 exploration. We then proposed our algorithm that explicitly targets uncertainty in reward differences
 479 most relevant for policy improvement. Under a multi-armed bandit setup of RLHF, we establish
 480 regret bounds of order $T^{(\beta+1)/(\beta+2)}$, which scales polynomially in all model parameters.

481 Our work opens several avenues for future investigation. An immediate question is whether the rate
 482 $T^{(\beta+1)/(\beta+2)}$ is minimax optimal, or if faster rates can be achieved. Another important direction is
 483 to refine the dependence on parameters such as A and r_{\max} , which may be improved with sharper
 484 analysis or alternative exploration schemes. Finally, our theoretical results are restricted to the bandit
 485 setting; extending the analysis to richer environments that incorporate a prompt space would be an
 exciting step toward bridging theory and practice in online RLHF.

486 REFERENCES
487

488 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for
489 reinforcement learning. In *International Conference on Machine Learning*, pp. 263–272. PMLR,
490 2017.

491 Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland,
492 Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning
493 from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp.
494 4447–4455. PMLR, 2024.

495 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
496 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless
497 assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*,
498 2022.

499 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method
500 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

501 Shicong Cen, Jincheng Mei, Katayoon Goshvadi, Hanjun Dai, Tong Yang, Sherry Yang, Dale
502 Schuurmans, Yuejie Chi, and Bo Dai. Value-incentivized preference optimization: A unified
503 approach to online and offline RLHF. In *The Thirteenth International Conference on Learning
504 Representations*, 2025. URL <https://openreview.net/forum?id=SQnitDuow6>.

505 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning
506 converts weak language models to strong language models. In *International Conference on
507 Machine Learning*, pp. 6621–6642. PMLR, 2024.

508 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
509 reinforcement learning from human preferences. *Advances in neural information processing
510 systems*, 30, 2017.

511 Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen
512 Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf.
513 *arXiv preprint arXiv:2405.07863*, 2024.

514 Yunzhen Feng, Ariel Kwiatkowski, Kunhao Zheng, Julia Kempe, and Yaqi Duan. Pilaf: Optimal
515 human preference sampling for reward modeling. *arXiv preprint arXiv:2502.04270*, 2025.

516 David A Freedman. On tail probabilities for martingales. *The Annals of Probability*, pp. 100–118,
517 1975.

518 Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre
519 Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from
520 online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.

521 Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep
522 Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. Camels in a changing climate:
523 Enhancing lm adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*, 2023.

524 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably
525 efficient? In *Advances in Neural Information Processing Systems*, pp. 4863–4873, 2018.

526 Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances
527 in applied mathematics*, 6(1):4–22, 1985.

528 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

529 Gen Li, Changxiao Cai, Yuxin Chen, Yuantao Gu, Yuting Wei, and Yuejie Chi. Is Q-learning
530 minimax optimal? a tight sample complexity analysis. *arXiv preprint arXiv:2102.06548*, 2021.

531 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

540 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
 541 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
 542 instructions with human feedback. *Advances in neural information processing systems*, 35:
 543 27730–27744, 2022.

544 Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From r to Q^* : Your language model
 545 is secretly a Q-function. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=kEVcNxtqXk>.

546 Corby Rosset, Ching-An Cheng, Arindam Mitra, Michael Santacroce, Ahmed Awadallah, and
 547 Tengyang Xie. Direct nash optimization: Teaching language models to self-improve with general
 548 preferences. *arXiv preprint arXiv:2404.03715*, 2024.

549 Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic
 550 exploration. *Advances in Neural Information Processing Systems*, 26, 2013.

551 Ruizhe Shi, Minhak Song, Runlong Zhou, Zihan Zhang, Maryam Fazel, and Simon S Du.
 552 Understanding the performance gap in preference learning: A dichotomy of rlhf and dpo. *arXiv*
 553 *preprint arXiv:2505.19770*, 2025.

554 Joel Tropp. Freedman’s inequality for matrix martingales. *Electronic Communications in
 555 Probability*, 16:262–270, 2011.

556 Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Hassan Awadallah,
 557 and Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q^* -
 558 approximation for sample-efficient RLHF. In *The Thirteenth International Conference
 559 on Learning Representations*, 2025. URL <https://openreview.net/forum?id=QYigQ6gXNw>.

560 Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang.
 561 Iterative preference learning from human feedback: Bridging theory and practice for rlhf under
 562 kl-constraint. *arXiv preprint arXiv:2312.11456*, 2023.

563 Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement
 564 learning without domain knowledge using value function bounds. In *International Conference on
 565 Machine Learning*, pp. 7304–7312. PMLR, 2019.

566 Shenao Zhang, Donghan Yu, Hiteshi Sharma, Han Zhong, Zhihan Liu, Ziyi Yang, Shuohang Wang,
 567 Hany Hassan Awadalla, and Zhaoran Wang. Self-exploring language models: Active preference
 568 elicitation for online alignment. *Transactions on Machine Learning Research*, 2025. ISSN 2835-
 569 8856. URL <https://openreview.net/forum?id=FoQK84nwY3>.

570 Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. Slic-hf:
 571 Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*, 2023.

572 Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, Karthik Ganesan, Wei-Lin Chiang, Jian Zhang,
 573 and Jiantao Jiao. Starling-7b: Improving helpfulness and harmlessness with rlaif. In *First
 574 Conference on Language Modeling*, 2024.

575 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
 576 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv
 577 preprint arXiv:1909.08593*, 2019.

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594 LLM USAGE
595

596 In preparing this paper, large language models (LLMs) were used as an assistive tool for minor
597 language polishing. All technical contributions, results, and conclusions are solely the work of the
598 authors.

600 A PROOF OF PROPOSITION 1
601

602 For each $t \geq 1$, define the event
603

$$604 \mathcal{E}_t := \{\text{no } a_0 \text{ is sampled in the first } t \text{ samples}\}.$$

605 We will show that for any $t \geq 1$,

$$606 \mathbb{P}(\mathcal{E}_t) \geq \frac{4}{9} (1 - \exp(-r_{\max}/\beta))^{2(t-1)}. \quad (\text{A.1})$$

607 Conditional on \mathcal{E}_t , it can be seen that $\ell(r, \mathcal{D}^{(t)})$ only depends on $r(a_1) - r(a_2)$. Now we study when
608 we fix $r(a_1) - r(a_2) \equiv \delta$ such that $\ell(r, \mathcal{D}^{(t)})$ is fixed, when is $J(\pi, r; \pi_{\text{cal}})$ maximized over both π
609 and r . By symmetry, we can assume without loss of generality that $\delta \geq 0$. We can compute
610

$$\begin{aligned} 611 J(\pi, r; \pi_{\text{cal}}) &= \mathbb{E}_{a \sim \pi}[r(a)] - \mathbb{E}_{a \sim \pi_{\text{cal}}}[r(a)] - \beta \text{KL}(\pi \parallel \pi_{\text{ref}}) \\ 612 &= [\pi(a_1) - p][r(a_1) - r(a_0)] + [\pi(a_2) - p][r(a_2) - r(a_0)] - \beta \text{KL}(\pi \parallel \pi_{\text{ref}}) \\ 613 &= [\pi(a_1) + \pi(a_2) - 2p][r(a_1) - r(a_0)] - \delta[\pi(a_2) - p] - \beta \text{KL}(\pi \parallel \pi_{\text{ref}}). \end{aligned}$$

614 For fixed π , we check which reward function r maximizes $J(\pi, r; \pi_{\text{cal}})$.

- 615 • When $\pi(a_1) + \pi(a_2) > 2p$, we know that

$$616 \max_r J(\pi, r; \pi_{\text{cal}}) = r_{\max}[\pi(a_1) + \pi(a_2) - 2p] - \delta[\pi(a_2) - p] - \beta \text{KL}(\pi \parallel \pi_{\text{ref}}), \quad (\text{A.2})$$

617 which is maximized at $r(a_1) = r_{\max}$, $r(a_2) = r_{\max} - \delta$ and $r(a_0) = 0$.

- 618 • When $\pi(a_1) + \pi(a_2) < 2p$, we know that

$$619 \max_r J(\pi, r; \pi_{\text{cal}}) = (r_{\max} - \delta)[2p - \pi(a_1) - \pi(a_2)] - \delta[\pi(a_2) - p] - \beta \text{KL}(\pi \parallel \pi_{\text{ref}}), \quad (\text{A.3})$$

620 which is maximized at $r(a_1) = \delta$, $r(a_2) = 0$ and $r(a_0) = r_{\max}$.

621 In addition, for any policy π such that $\pi(a_1) + \pi(a_2) < 2p$, by considering another policy π' defined
622 as $\pi'(a_1) = 2p - \pi(a_2)$ and $\pi'(a_2) = 2p - \pi(a_1)$, we have

$$\begin{aligned} 623 \max_r J(\pi', r; \pi_{\text{cal}}) - \max_r J(\pi, r; \pi_{\text{cal}}) \\ 624 &= r_{\max}[\pi'(a_1) + \pi'(a_2) - 2p] - \delta[\pi'(a_2) - p] - \beta \text{KL}(\pi' \parallel \pi_{\text{ref}}) \\ 625 &\quad - (r_{\max} - \delta)[2p - \pi(a_1) - \pi(a_2)] + \delta[\pi(a_2) - p] + \beta \text{KL}(\pi \parallel \pi_{\text{ref}}) \\ 626 &= \beta[\text{KL}(\pi \parallel \pi_{\text{ref}}) - \text{KL}(\pi' \parallel \pi_{\text{ref}})]. \end{aligned}$$

627 Here the first relation follows from (A.2), (A.3) and the fact that $\pi'(a_1) + \pi'(a_2) > 2p$. Let $x =$
628 $\pi(a_1)$ and $y = \pi(a_2)$. Let

$$\begin{aligned} 629 f(x, y) &:= \text{KL}(\pi \parallel \pi_{\text{ref}}) - \text{KL}(\pi' \parallel \pi_{\text{ref}}) \\ 630 &= x \log x + y \log y + (1 - x - y) \log(1 - x - y) - (2p - x) \log(2p - x) \\ 631 &\quad - (2p - y) \log(2p - y) - (1 - 4p + x + y) \log(1 + x + y - 4p). \end{aligned}$$

632 By elementary analysis, it is straightforward to check that $f(x, y) > 0$ for any $x, y > 0$ satisfying
633 $x + y < 2p$. Therefore we have

$$634 \max_r J(\pi', r; \pi_{\text{cal}}) > \max_r J(\pi, r; \pi_{\text{cal}}).$$

635 Therefore in order to maximize $\ell(r, \mathcal{D}^{(t)}) + \alpha J^*(r; \pi_{\text{cal}})$, the following statement always holds
636 regardless of the value of δ :

$$637 r^{(t+1)}(a_0) = 0, \quad \max \{r^{(t+1)}(a_1), r^{(t+1)}(a_2)\} = r_{\max}.$$

648 This immediately implies that
 649

$$650 \pi^{(t+1)}(a_0) = \frac{\exp(r^{(t+1)}(a_0)/\beta)}{\exp(r^{(t+1)}(a_0)/\beta) + \exp(r^{(t+1)}(a_1)/\beta) + \exp(r^{(t+1)}(a_2)/\beta)} \leq \frac{1}{2 + \exp(r_{\max}/\beta)}.$$

652 Therefore conditional on \mathcal{E}_t , we know that
 653

$$654 \mathbb{P}(\mathcal{E}_{t+1}|\mathcal{E}_t) \geq (1 - \pi^{(t+1)}(a_0))^2 \geq \left(\frac{1}{1 + \exp(-r_{\max}/\beta)} \right)^2 \geq (1 - \exp(-r_{\max}/\beta))^2.$$

655 This relation, together with
 656

$$657 \mathbb{P}(\mathcal{E}_0) = (\pi^{(1)}(a_1) + \pi^{(1)}(a_2))^2 = \frac{4}{9},$$

659 establishes the statement (A.1). This immediately implies that, for any $t \leq \exp(r_{\max}/\beta)/2$,
 660

$$663 \mathbb{P}(\mathcal{E}_t) \geq \frac{4}{9} (1 - \exp(-r_{\max}/\beta))^{2(t-1)} \geq \frac{4}{9} (1 - \exp(-r_{\max}/\beta))^{\exp(r_{\max}/\beta)} \geq \frac{4}{9e} \geq 0.16.$$

666 Finally, when \mathcal{E}_t holds, we have
 667

$$668 J(\pi^*; r^*, \pi_{\text{cal}}) - J(\pi^{(t)}; r^*, \pi_{\text{cal}}) = \pi^*(a_0) - \pi^{(t)}(a_0) - \beta \text{KL}(\pi^* \parallel \pi_{\text{ref}}) + \beta \text{KL}(\pi^{(t)} \parallel \pi_{\text{ref}}).$$

669 We have
 670

$$671 \pi^*(a_0) = \frac{\exp(1/\beta)}{\exp(1/\beta) + 2}, \quad \pi^*(a_1) = \pi^*(a_2) = \frac{1}{\exp(1/\beta) + 2}.$$

673 Therefore we have
 674

$$675 \text{KL}(\pi^* \parallel \pi_{\text{ref}}) = \log 3 + \pi^*(a_0) \log \pi^*(a_0) + \pi^*(a_1) \log \pi^*(a_1) + \pi^*(a_2) \log \pi^*(a_2) \\ 676 = \log 3 + \frac{\exp(1/\beta)}{\exp(1/\beta) + 2} \log \frac{\exp(1/\beta)}{\exp(1/\beta) + 2} + \frac{2}{\exp(1/\beta) + 2} \log \frac{1}{\exp(1/\beta) + 2} \\ 677 = \log 3 + \beta^{-1} \frac{\exp(1/\beta)}{\exp(1/\beta) + 2} - \log[\exp(1/\beta) + 2].$$

679 In addition, when \mathcal{E}^{t-1} happens, we know that
 680

$$682 \text{KL}(\pi^{(t)} \parallel \pi_{\text{ref}}) = \log 3 + \pi^{(t)}(a_0) \log \pi^{(t)}(a_0) + \pi^{(t)}(a_1) \log \pi^{(t)}(a_1) + \pi^{(t)}(a_2) \log \pi^{(t)}(a_2) \\ 683 \stackrel{(i)}{\geq} \log 3 + \pi^{(t)}(a_0) \log \pi^{(t)}(a_0) + [\pi^{(t)}(a_1) + \pi^{(t)}(a_2)] \log \frac{\pi^{(t)}(a_1) + \pi^{(t)}(a_2)}{2} \\ 684 = \log 3 + \pi^{(t)}(a_0) \log \pi^{(t)}(a_0) + [1 - \pi^{(t)}(a_0)] \log \frac{1 - \pi^{(t)}(a_0)}{2} \\ 685 \stackrel{(ii)}{\geq} \log 3 - \log 2 - 0.16.$$

687 Here step (i) uses Jensen's inequality for convex function $f(x) = x \log x$; step (ii) holds since the
 688 function $g(x) = x \log x + (1 - x) \log(1 - x)/2$ is monotonically decreasing for $0 < x < 1/3$, and
 689 we have
 690

$$691 \pi^{(t)}(a_0) \leq \frac{1}{2 + \exp(r_{\max}/\beta)} \leq \frac{1}{2 + \exp(3)} \leq 0.046$$

695 provided that $r_{\max}/\beta \geq 3$. We have
 696

$$697 J(\pi^*; r^*, \pi_{\text{cal}}) - J(\pi^{(t)}; r^*, \pi_{\text{cal}}) = \pi^*(a_0) - \pi^{(t)}(a_0) - \beta \text{KL}(\pi^* \parallel \pi_{\text{ref}}) + \beta \text{KL}(\pi^{(t)} \parallel \pi_{\text{ref}}) \\ 698 \geq \beta \log(\exp(1/\beta) + 2) - (\log 2 + 0.16)\beta - 0.046 \\ 699 \geq 1/2,$$

701 where the last relation holds for any $\beta > 0$.

702 **B PROOF OF PROPOSITION 2**
 703

704 Let $T = \min\{\kappa, \exp(r_{\max})/2\}$, and define the events
 705

706
$$\mathcal{A} := \{a_2^t = a_0 \text{ for all } 1 \leq t \leq T\}$$

 707

and

709
$$\mathcal{E} := \{a_1^t \succ a_2^t \text{ or } a_1^t = a_2^t \text{ for all } 1 \leq t \leq T\}.$$

 710

We can check that when $\kappa \geq 5$,

711
$$\mathbb{P}(\mathcal{A}) = [\pi_{\text{ref}}(a_0)]^T \leq (1 - 2\kappa^{-1})^\kappa \geq \frac{1}{16}.$$

 712

713 Conditional on \mathcal{A} , we know that when $r_{\max} \geq 1$,
 714

715
$$\mathbb{P}(\mathcal{E} | \mathcal{A}) \geq \left(\frac{\exp(r_{\max} - 1)}{1 + \exp(r_{\max} - 1)} \right)^T \geq \left(\frac{\exp(r_{\max} - 1)}{1 + \exp(r_{\max} - 1)} \right)^{\exp(r_{\max})/2} \geq \frac{1}{4}.$$

 716

717 Conditional on \mathcal{A} and \mathcal{E}_t , for any $1 \leq t \leq T-1$, all the preference data in $\mathcal{D}^{(t)}$ are of form $a_1^t \succ a_2^t$.
 718 In this case, it is straightforward to check that the reward function that maximizes $\ell(r, \mathcal{D}^{(t)}) +$
 719 $\alpha J^*(r; \pi_{\text{ref}})$ is
 720

721
$$r^{(t+1)}(a_0) = 0, \quad r^{(t+1)}(a_1) = r^{(t+1)}(a_2) = r_{\max}.$$

 722

This immediately implies that

723
$$\pi^{(t+1)}(a_0) = \frac{\kappa - 2}{\kappa - 2 + 2 \exp(r_{\max}/\beta)}, \quad \pi^{(t+1)}(a_1) = \pi^{(t+1)}(a_2) = \frac{\exp(r_{\max}/\beta)}{\kappa - 2 + 2 \exp(r_{\max}/\beta)}.$$

 724

On the other hand, we know that

725
$$\begin{aligned} \pi^*(a_0) &= \frac{\kappa - 2}{\kappa - 2 + \exp(r_{\max}/\beta) + \exp((r_{\max} - 2)/\beta)}, \\ 726 \pi^*(a_1) &= \frac{\exp(r_{\max}/\beta)}{\kappa - 2 + \exp(r_{\max}/\beta) + \exp((r_{\max} - 2)/\beta)}, \\ 727 \pi^*(a_2) &= \frac{\exp((r_{\max} - 2)/\beta)}{\kappa - 2 + \exp(r_{\max}/\beta) + \exp((r_{\max} - 2)/\beta)}. \end{aligned}$$

 728

729 For any $2 \leq t \leq T$, we first lower bound
 730

731
$$J(\pi^*; r^*, \pi_{\text{ref}}) - J(\pi^{(t)}; r^*, \pi_{\text{ref}}) \geq J(\pi_{\theta^*}; r^*, \pi_{\text{ref}}) - J(\pi_1; r^*, \pi_{\text{ref}}) \quad (\text{B.1})$$

 732

733 for any $\theta^* \in [0, 1]$, where we define $\pi_\theta := \theta\pi^{(t)} + (1 - \theta)\pi^*$, and the above relation follows from
 734 the optimality of π^* . Recall the definition
 735

736
$$J(\pi; r^*, \pi_{\text{ref}}) = \pi(a_1)r_{\max} + \pi(a_2)(r_{\max} - 2) - \beta \sum_{i=0}^2 \pi(a_i) \log \frac{\pi(a_i)}{\pi_{\text{ref}}(a_i)},$$

 737

738 we can compute
 739

740
$$\nabla_\pi J(\pi; r^*, \pi_{\text{ref}}) = \begin{bmatrix} r^*(a_0) - \beta \log[\pi(a_0)/\pi_{\text{ref}}(a_0)] - \beta \\ r^*(a_1) - \beta \log[\pi(a_1)/\pi_{\text{ref}}(a_1)] - \beta \\ r^*(a_2) - \beta \log[\pi(a_2)/\pi_{\text{ref}}(a_2)] - \beta \end{bmatrix}$$

 741

742 and

743
$$\nabla_\pi^2 J(\pi; r^*, \pi_{\text{ref}}) = -\beta \text{diag}\{\pi(a_0), \pi(a_1), \pi(a_2)\}^{-1}.$$

 744

745 It is straightforward to check that
 746

747
$$\nabla_\pi J(\pi^{(t)}; r^*, \pi_{\text{ref}}) = \begin{bmatrix} 0 \\ 0 \\ -1 \end{bmatrix} + \text{const} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}. \quad (\text{B.2})$$

 748

756 Since $\pi^{(t)}(a_0) < \pi^*(a_0)$, $\pi^{(t)}(a_1) < \pi^*(a_1)$ and $\pi^{(t)}(a_2) > \pi^*(a_2)$, we know that for any $\theta \in$
 757 $[0, \theta^*]$

$$759 \nabla_\pi^2 J(\pi_\theta; r^*, \pi_{\text{ref}}) \succeq -\beta \text{diag}\{\pi^{(t)}(a_0), \pi^{(t)}(a_1), \theta^* \pi^{(t)}(a_2) + (1 - \theta^*) \pi^*(a_2)\}^{-1}. \quad (\text{B.3})$$

760 Therefore we have

$$\begin{aligned} 762 J(\pi_{\theta^*}; r^*, \pi_{\text{ref}}) - J(\pi_1; r^*, \pi_{\text{ref}}) &\stackrel{\text{(i)}}{\geq} \theta^* \nabla_\pi J(\pi^{(t)}; r^*, \pi_{\text{ref}})^\top (\pi^* - \pi^{(t)}) \\ 763 &\quad - \frac{\beta \theta^{*2}}{2} (\pi^* - \pi^{(t)})^\top \text{diag}\{\pi^{(t)}(a_0), \pi^{(t)}(a_1), \theta^* \pi^{(t)}(a_2) + (1 - \theta^*) \pi^*(a_2)\}^{-1} (\pi^* - \pi^{(t)}) \\ 764 &\stackrel{\text{(ii)}}{\geq} \theta^* [\pi^{(t)}(a_2) - \pi^*(a_2)] - \frac{\beta \theta^{*2}}{2} [\pi^*(a_0) + \pi^*(a_1) + \frac{9}{16} \pi^{(t)}(a_2) / \theta^*] \\ 765 &= \theta^* [\pi^{(t)}(a_2) - \pi^*(a_2)] - \frac{9}{32} \beta \theta^* \pi^{(t)}(a_2) - \frac{\beta \theta^{*2}}{2} [1 - \pi^*(a_2)] \\ 766 &= \left(1 - \frac{9}{32} \beta\right) \theta^* \pi^{(t)}(a_2) - \left(1 - \frac{\beta \theta^*}{2}\right) \theta^* \pi^*(a_2) - \frac{\beta \theta^{*2}}{2} \\ 767 &\stackrel{\text{(iii)}}{\geq} \left(\frac{3}{4} - \frac{9}{32} \beta + \frac{\beta \theta^*}{8}\right) \theta^* \pi^{(t)}(a_2) - \frac{\beta \theta^{*2}}{2} \stackrel{\text{(iv)}}{\geq} \frac{15}{32} \theta^* \pi^{(t)}(a_2) - \frac{\theta^{*2}}{2}. \end{aligned} \quad (\text{B.4})$$

775 Here step (i) follows from the Taylor expansion and (B.3); step (ii) utilizes (B.2) and as well as the
 776 following relations

$$777 \pi^{(t)}(a_0) \leq \pi^*(a_0) \leq 2\pi^{(t)}(a_0), \quad \pi^{(t)}(a_1) \leq \pi^*(a_1) \leq 2\pi^{(t)}(a_1)$$

779 and when $\beta \leq 1$,

$$780 \pi^*(a_2) \leq \frac{2}{\exp(2/\beta) + 1} \pi^{(t)}(a_2) \leq \frac{1}{4} \pi^{(t)}(a_2); \quad (\text{B.5})$$

782 steps (iii) and (iv) follows from (B.5) and $\beta \leq 1$. When $\kappa \leq \exp(r_{\max} \beta)$, we have

$$784 \pi^{(t)}(a_2) = \frac{\exp(r_{\max}/\beta)}{\kappa - 2 + 2 \exp(r_{\max}/\beta)} \geq \frac{\exp(r_{\max}/\beta)}{3 \exp(r_{\max}/\beta) - 2} \geq \frac{1}{3}. \quad (\text{B.6})$$

786 By taking (B.1), (B.4) and (B.6) collectively, we have

$$788 (\pi^*; r^*, \pi_{\text{ref}}) - J(\pi^{(t)}; r^*, \pi_{\text{ref}}) \geq \frac{5}{32} \theta^* - \frac{\theta^{*2}}{2} \geq \frac{25}{2048} > 0.01$$

790 where we take $\theta^* = 5/32$.

792 C PROOF OF AUXILIARY LEMMAS

795 C.1 PROOF OF LEMMA 1

796 We first express

$$798 X_i := \log \frac{\sigma(r^*(a_+^i) - r^*(a_-^i))}{\sigma(r(a_+^i) - r(a_-^i))} = \mathbb{1}\{a_1^i \succ a_2^i\} \log \frac{\sigma(r^*(a_1^i) - r^*(a_2^i))}{\sigma(r(a_1^i) - r(a_2^i))} + \mathbb{1}\{a_1^i \prec a_2^i\} \log \frac{\sigma(r^*(a_2^i) - r^*(a_1^i))}{\sigma(r(a_2^i) - r(a_1^i))}.$$

801 It is straightforward to check that

$$\begin{aligned} 802 \mathbb{E}[X_i | a_1^i, a_2^i] &= \mathbb{P}(a_1^i \succ a_2^i | a_1^i, a_2^i) \log \frac{\sigma(r^*(a_1^i) - r^*(a_2^i))}{\sigma(r(a_1^i) - r(a_2^i))} + \mathbb{P}(a_1^i \prec a_2^i | a_1^i, a_2^i) \log \frac{\sigma(r^*(a_2^i) - r^*(a_1^i))}{\sigma(r(a_2^i) - r(a_1^i))} \\ 803 &= \sigma(r^*(a_1^i) - r^*(a_2^i)) \log \frac{\sigma(r^*(a_1^i) - r^*(a_2^i))}{\sigma(r(a_1^i) - r(a_2^i))} + \sigma(r^*(a_2^i) - r^*(a_1^i)) \log \frac{\sigma(r^*(a_2^i) - r^*(a_1^i))}{\sigma(r(a_2^i) - r(a_1^i))} \\ 804 &= \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))). \end{aligned}$$

808 and

$$809 |X_i| \leq |\log(1 + \exp(-r(a_+^i) + r(a_-^i)))| \leq 2r_{\max}.$$

810 In addition, we can compute the variance
 811

$$\begin{aligned}
 812 \text{Var}(X_i | a_1^i, a_2^i) &= \sigma(r^*(a_1^i) - r^*(a_2^i))\sigma(r^*(a_2^i) - r^*(a_1^i)) \left[\log \frac{\sigma(r^*(a_1^i) - r^*(a_2^i))}{\sigma(r(a_1^i) - r(a_2^i))} - \log \frac{\sigma(r^*(a_2^i) - r^*(a_1^i))}{\sigma(r(a_2^i) - r(a_1^i))} \right]^2 \\
 813 &= \sigma(r^*(a_1^i) - r^*(a_2^i))\sigma(r^*(a_2^i) - r^*(a_1^i)) \left[\log \frac{\sigma(r^*(a_1^i) - r^*(a_2^i))}{\sigma(r^*(a_2^i) - r^*(a_1^i))} - \log \frac{\sigma(r(a_1^i) - r(a_2^i))}{\sigma(r(a_2^i) - r(a_1^i))} \right]^2 \\
 814 &= \sigma(r^*(a_1^i) - r^*(a_2^i))\sigma(r^*(a_2^i) - r^*(a_1^i)) [r(a_1^i) - r(a_2^i) - r^*(a_1^i) + r^*(a_2^i)]^2.
 \end{aligned}$$

815 In view of Lemma 5, we have
 816

$$\begin{aligned}
 817 \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) \\
 818 &\geq \frac{1}{4}\sigma(r^*(a_1^i) - r^*(a_2^i))\sigma(r^*(a_2^i) - r^*(a_1^i)) \\
 819 &\quad \cdot \min \left\{ |r(a_1^i) - r(a_2^i) - r^*(a_1^i) + r^*(a_2^i)|, [r(a_1^i) - r(a_2^i) - r^*(a_1^i) + r^*(a_2^i)]^2 \right\} \\
 820 &\geq \frac{1}{16r_{\max}}\sigma(r^*(a_1^i) - r^*(a_2^i))\sigma(r^*(a_2^i) - r^*(a_1^i)) [r(a_1^i) - r(a_2^i) - r^*(a_1^i) + r^*(a_2^i)]^2,
 \end{aligned} \tag{C.1}$$

821 where the last step follows from $|r(a_1^i) - r(a_2^i) - r^*(a_1^i) + r^*(a_2^i)| \leq 4r_{\max}$. Therefore we have
 822

$$\text{Var}(X_i | a_1^i, a_2^i) \leq 16r_{\max} \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))).$$

823 In addition, we have the following deterministic bound
 824

$$\sum_{i=1}^t \text{Var}(X_i | a_1^i, a_2^i) \leq 16tr_{\max}^2.$$

825 By the Freedman's inequality (cf. Lemma 6), for any fixed r , with probability exceeding $1 - \delta$,
 826

$$\begin{aligned}
 827 |\Delta_t(r)| &\leq \left| \sum_{i=1}^t (X_i - \mathbb{E}[X_i | a_1^i, a_2^i]) \right| \\
 828 &\leq C_2 \sqrt{\sum_{i=1}^t r_{\max} \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) \log \frac{\log t}{\delta} + C_2 r_{\max} \log \frac{\log t}{\delta}}
 \end{aligned}$$

829 for some sufficiently large constant $C_2 > 0$.
 830

831 C.2 PROOF OF LEMMA 2

832 For any fixed $r : \mathcal{A} \rightarrow [\pm r_{\max}]$, with probability exceeding $1 - \delta$ we have
 833

$$\begin{aligned}
 834 |\Delta_t(r)| &\stackrel{(i)}{\leq} C_2 \sqrt{\sum_{i=1}^t r_{\max} \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) \log \frac{\log T}{\delta} + C_2 r_{\max} \log \frac{\log T}{\delta}} \\
 835 &\stackrel{(ii)}{\leq} \frac{1}{2} \sum_i \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) + 2C_2^2 r_{\max} \log \frac{\log T}{\delta}.
 \end{aligned}$$

836 Here step (i) follows from Lemma 1, and step (ii) utilizes the AM-GM inequality. This immediately
 837 implies that
 838

$$\begin{aligned}
 839 \ell(r, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)}) &= \sum_{i=1}^t \log \frac{\sigma(r(a_+^i) - r(a_-^i))}{\sigma(r^*(a_+^i) - r^*(a_-^i))} \\
 840 &= - \sum_{i=1}^t \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) - \Delta_t(r)
 \end{aligned}$$

$$\begin{aligned}
& \leq -\frac{1}{2} \sum_{i=1}^t \mathsf{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) + 2C_2^2 r_{\max} \log \frac{\log T}{\delta}.
\end{aligned} \tag{C.2}$$

Then we explore the Lipschitzness continuity of the above functionals of r . For any two fixed reward functions $r, r' : \mathcal{A} \rightarrow [\pm r_{\max}]$, we have

$$\begin{aligned}
|\ell(r, \mathcal{D}^{(t)}) - \ell(r', \mathcal{D}^{(t)})| &= \sum_{i=1}^t |\log[\sigma(r(a_+^i) - r(a_-^i))] - \log[\sigma(r'(a_+^i) - r'(a_-^i))]| \\
&\leq \sum_{i=1}^t |r(a_+^i) - r(a_-^i) - r'(a_+^i) + r'(a_-^i)| \leq 2T \|r - r'\|_{\infty},
\end{aligned} \tag{C.3}$$

where the penultimate step follows from $d \log(\sigma(x))/dx = \sigma(-x) \leq 1$. Similarly, for any $x, y, \delta \in \mathbb{R}$, we have

$$\begin{aligned}
|\mathsf{KL}(\sigma(x) \parallel \sigma(y)) - \mathsf{KL}(\sigma(x) \parallel \sigma(y + \delta))| &= \left| \sigma(x) \log \frac{\sigma(y + \delta)}{\sigma(y)} + (1 - \sigma(x)) \log \frac{1 - \sigma(y + \delta)}{1 - \sigma(y)} \right| \\
&\leq \sigma(x)|\delta| + (1 - \sigma(x))|\delta| = |\delta|.
\end{aligned}$$

This implies that

$$\begin{aligned}
&\left| \sum_{i=1}^t \mathsf{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) \right. \\
&\quad \left. - \sum_{i=1}^t \mathsf{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r'(a_1^i) - r'(a_2^i))) \right| \leq 2 \|r - r'\|_{\infty}.
\end{aligned} \tag{C.4}$$

Let $\mathcal{N}_{\varepsilon}$ be an ε -net of $[-r_{\max}, r_{\max}]^A$ (or equivalently, the function space of $r : \mathcal{A} \rightarrow [\pm r_{\max}]$) under the ℓ_{∞} norm such that $|\mathcal{N}_{\varepsilon}| \leq (2r_{\max}/\varepsilon)^A$. By standard union bound argument and (C.2), with probability exceeding $1 - \delta$,

$$\ell(r, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)}) \leq -\frac{1}{2} \sum_{i=1}^t \mathsf{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) + 2C_2^2 r_{\max} \log \frac{|\mathcal{N}_{\varepsilon}| \log T}{\delta} \tag{C.5}$$

holds for any $r \in \mathcal{N}_{\varepsilon}$. This implies that for any $r : \mathcal{A} \rightarrow [\pm r_{\max}]$, there exists $r_0 \in \mathcal{N}_{\varepsilon}$ such that $\|r - r'\| \leq \varepsilon$, hence

$$\begin{aligned}
\ell(r, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)}) &\stackrel{(i)}{\leq} \ell(r_0, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)}) + 2T\varepsilon \\
&\stackrel{(ii)}{\leq} -\frac{1}{2} \sum_{i=1}^t \mathsf{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r_0(a_1^i) - r_0(a_2^i))) + 2C_2^2 r_{\max} \log \frac{|\mathcal{N}_{\varepsilon}| \log T}{\delta} + 2T\varepsilon \\
&\stackrel{(iii)}{\leq} -\frac{1}{2} \sum_{i=1}^t \mathsf{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) + 2C_2^2 r_{\max} \log \frac{|\mathcal{N}_{\varepsilon}| \log T}{\delta} + 4T\varepsilon \\
&\stackrel{(iv)}{\leq} -\frac{1}{2} \sum_{i=1}^t \mathsf{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r(a_1^i) - r(a_2^i))) + C_3 A r_{\max} \log T.
\end{aligned}$$

Here step (i) utilizes (C.3); step (ii) follows from $r_0 \in \mathcal{N}_{\varepsilon}$ and the uniform concentration bound (C.5); step (iii) uses (C.4); step (iv) holds as long as $C_3 \gg 2C_2^2$, where we let $\varepsilon = A r_{\max}/T$ and $\delta = T^{-10}$. This completes the proof.

C.3 PROOF OF LEMMA 3

When $N_t(a_+, a_-) \geq 100C_4 \alpha_t r_{\max}$, we have

$$\mathsf{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t)}(a_+) - r^{(t)}(a_-))) \leq \frac{1}{100}. \tag{C.6}$$

918 Now we assert that $r^{(t)}(a_-) - r^{(t)}(a_+) < 0.5$ for any $t \geq t_0$. This is because, if $r^{(t)}(a_-) - r^{(t)}(a_+) \geq 0.5$, we have

$$921 \text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t)}(a_+) - r^{(t)}(a_-))) = \text{KL}(\sigma(r^*(a_-) - r^*(a_+)) \parallel \sigma(r^{(t)}(a_-) - r^{(t)}(a_+))) \\ 922 \geq \text{KL}(\sigma(0) \parallel \sigma(0.5)) > \frac{1}{100}.$$

923 Here we use the fact that $r^*(a_-) - r^*(a_+) \leq 0$. This contradicts with (C.6). Hence we have

$$926 r^{(t)}(a_-) - r^{(t)}(a_+) < 0.5. \quad (C.7)$$

928 Let $p := \sigma(r^*(a_-) - r^*(a_+))$ and $q := \sigma(r^{(t)}(a_-) - r^{(t)}(a_+))$. We have

$$930 \exp(r^{(t)}(a_-) - r^{(t)}(a_+)) \stackrel{(i)}{\leq} 3\sigma(r^{(t)}(a_-) - r^{(t)}(a_+)) = 3q \stackrel{(ii)}{\leq} 6p + \text{KL}(p \parallel q) \\ 931 = 6\sigma(r^*(a_-) - r^*(a_+)) + 24\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t)}(a_+) - r^{(t)}(a_-))).$$

933 Here step (i) follows from (C.7), while step (ii) holds trivially when $q \leq 2p$, and when $q > 2p$ we
934 have

$$935 \text{KL}(p \parallel q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \geq \frac{(q-p)^2}{2q} \geq \frac{1}{8}q.$$

938 Finally, for any $t_0 \leq t_1 < t_2 \leq T$, we can upper bound

$$940 N_{t_2}(a_+, a_-) - N_{t_1}(a_+, a_-) \leq \sum_{i=t_1+1}^{t_2} X_i \quad \text{where} \quad X_i := \mathbb{1}\{\text{a}_- \text{ is sampled in the } i\text{-th iteration}\}.$$

943 It is straightforward to check that $X_i - \mathbb{E}[X_i | \mathcal{F}_{i-1}]$ is a martingale difference sequence, and by the
944 Azuma-Hoeffding inequality, with probability exceeding $1 - O(T^{-100})$ we have

$$945 \sum_{i=t_1+1}^{t_2} (X_i - \mathbb{E}[X_i | \pi^{(i)}, \pi^{(i-1)}]) \leq \tilde{C} \sqrt{T \log T}$$

948 for some universal constant $\tilde{C} > 0$. In addition, we have

$$950 \mathbb{E}[X_i | \pi^{(i)}, \pi^{(i-1)}] \leq \frac{\pi^{(i)}(a_-)}{\pi^{(i)}(a_-) + \pi^{(i)}(a_+)} + \frac{\pi^{(i-1)}(a_-)}{\pi^{(i-1)}(a_-) + \pi^{(i-1)}(a_+)}.$$

953 For each $t \in [T]$, we have

$$955 \frac{\pi^{(t)}(a_-)}{\pi^{(t)}(a_-) + \pi^{(t)}(a_+)} \stackrel{(i)}{=} \frac{\pi_{\text{ref}}(a_-) \exp(r^{(t)}(a_-)/\beta)}{\pi_{\text{ref}}(a_-) \exp(r^{(t)}(a_-)/\beta) + \pi_{\text{ref}}(a_+) \exp(r^{(t)}(a_+)/\beta)} \\ 956 \leq \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \exp\left((r^{(t)}(a_-) - r^{(t)}(a_+))/\beta\right) \\ 957 \leq \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \left[6\sigma(r^*(a_-) - r^*(a_+)) + 24\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t)}(a_+) - r^{(t)}(a_-)))\right]^{1/\beta}.$$

962 Here step (i) utilizes (2.4), while step (ii) follows from (C.8). Hence we have

$$964 N_{t_2}(a_+, a_-) - N_{t_1}(a_+, a_-) \leq 2 \sum_{t=t_1+1}^{t_2} \frac{\pi^{(t)}(a_-)}{\pi^{(t)}(a_-) + \pi^{(t)}(a_+)} + 2\tilde{C} \sqrt{T \log T} \\ 965 \leq C_5^{1/\beta} \sum_{t=t_1+1}^{t_2} \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \left[\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t)}(a_+) - r^{(t)}(a_-))) \right]^{\frac{1}{\beta}} \\ 966 + \sigma(r^*(a_-) - r^*(a_+))^{\frac{1}{\beta}} + C_5 \sqrt{T \log T}$$

971 for some sufficiently large constant $C_5 > 0$.

972 C.4 PROOF OF LEMMA 4
973974 Let t_0 be the first iteration such that
975

976
$$N_{t_0}(a_+, a_-) \geq \min \left\{ \frac{1}{2} N_T(a_+, a_-), 100C_4\alpha_T r_{\max} \right\}. \quad (\text{C.9})$$

977

978 In what follows, we establish the desired result under two different cases: $N_T(a_+, a_-)$ being larger
979 or smaller than $c_0 \exp(r^*(a_+) - r^*(a_-))\alpha_T r_{\max}$ for some sufficiently large constant $c_0 > 0$.
980981 **Case 1.** When $N_T(a_+, a_-) \leq c_0 \exp(r^*(a_+) - r^*(a_-))\alpha_T r_{\max}$, it is straightforward to show that
982

983
$$\zeta(a_+, a_-) \leq N_T(a_+, a_-)r_{\max} \leq C_6 \exp(r^*(a_+) - r^*(a_-))\alpha_T r_{\max}^2 = C_6 \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T r_{\max}^2. \quad (\text{C.10})$$

984

985 In addition, we have
986

987
$$\sigma(r^*(a_-) - r^*(a_+)) \leq \exp(r^*(a_-) - r^*(a_+)) \leq \frac{c_0 \alpha_T r_{\max}}{N_T(a_+, a_-)}.$$

988

989 In addition, for any $t_0 \leq t \leq T$, we can use (5.8) to show that
990

991
$$\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \| \sigma(r^{(t)}(a_+) - r^{(t)}(a_-))) \leq \frac{C_4 \alpha_T r_{\max}}{N_t(a_+, a_-)} \leq \frac{2C_4 \alpha_T r_{\max}}{N_T(a_+, a_-)}. \quad (\text{C.11})$$

992

993 By taking $t_1 = t_0 - 1$ and $t_2 = T$ in Lemma 3, we have
994

995
$$\begin{aligned} N_T(a_+, a_-) &\stackrel{\text{(i)}}{\leq} 2[N_T(a_+, a_-) - N_{t_0-1}(a_+, a_-)] \\ &\stackrel{\text{(ii)}}{\leq} 4T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \left(\frac{C_5 \max\{c_0, 2C_4\} \alpha_T r_{\max}}{N_T(a_+, a_-)} \right)^{1/\beta} + 2C_5 \sqrt{T \log T}. \end{aligned}$$

996

997 Here step (i) follows from the definition of t_0 (cf. (C.9)), while step (ii) uses the above two bounds
998 and Lemma 3 with $t_1 = t_0 - 1$ and $t_2 = T$. This immediately implies that
999

1000
$$N_T(a_+, a_-) \leq C_7 \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} (\alpha_T r_{\max})^{\frac{1}{\beta+1}} + C_7 \sqrt{T \log T}$$

1001

1002 for some sufficiently large constant $C_7 > 0$. This leads to
1003

1004
$$\zeta(a_+, a_-) \leq N_T(a_+, a_-)r_{\max} \leq C_7 \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{\beta+2}{\beta+1}} + C_7 \sqrt{T \log T} r_{\max}. \quad (\text{C.12})$$

1005

1006 **Case 2.** When $N_T(a_+, a_-) > c_0 \exp(r^*(a_+) - r^*(a_-))\alpha_T r_{\max}$, we have
1007

1008
$$\begin{aligned} \exp(r^*(a_+) - r^*(a_-))\alpha_T r_{\max} &\leq \frac{1}{c_0} N_T(a_+, a_-) \stackrel{\text{(i)}}{\leq} \frac{2}{c_0} [N_T(a_+, a_-) - N_{t_0-1}(a_+, a_-)] \\ &\stackrel{\text{(ii)}}{\leq} \frac{2C_5^{1/\beta}}{c_0} T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \left[\sigma(r^*(a_-) - r^*(a_+))^{1/\beta} + \left(\frac{2C_4 \alpha_T r_{\max}}{N_T(a_+, a_-)} \right)^{1/\beta} \right] + \frac{2C_5}{c_0} \sqrt{T \log T} \\ &\stackrel{\text{(iii)}}{\leq} \frac{4}{c_0} \max\{C_5, 2C_4 C_5/c_0\}^{1/\beta} T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \exp(r^*(a_-) - r^*(a_+))^{1/\beta} + \frac{2C_5}{c_0} \sqrt{T \log T}. \end{aligned}$$

1009

1010 Here step (i) follows from the definition of t_0 (cf. (C.9)); step (ii) utilizes Lemma 3 with $t_1 = t_0 - 1$
1011 and $t_2 = T$, as well as (C.11); step (iii) holds since $\sigma(r^*(a_-) - r^*(a_+)) \leq \exp(r^*(a_-) - r^*(a_+))$
1012 and
1013

1014
$$\frac{2C_4 \alpha_T r_{\max}}{N_T(a_+, a_-)} \leq \frac{2C_4 \alpha_T r_{\max}}{c_0 \exp(r^*(a_+) - r^*(a_-))\alpha_T r_{\max}} \leq \frac{2C_4}{c_0} \exp(r^*(a_-) - r^*(a_+)).$$

1015

1016 This immediately implies that for some sufficiently large constant $C_8 > 0$, we have
1017

1018
$$\frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} = \exp(r^*(a_+) - r^*(a_-)) \leq C_8 \left(\frac{\pi_{\text{ref}}(a_-)T}{\pi_{\text{ref}}(a_+) \alpha_T r_{\max}} \right)^{\frac{\beta}{\beta+1}} + C_8 \frac{\sqrt{T \log T}}{\alpha_T r_{\max}}. \quad (\text{C.13})$$

1019

1026 Similar to (C.1), we can show that
1027
1028 $\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r(a_+) - r(a_-))) = \text{KL}(\sigma(r^*(a_-) - r^*(a_+)) \parallel \sigma(r(a_-) - r(a_+)))$
1029 $\stackrel{(a)}{\geq} \frac{1}{16r_{\max}} \sigma(r^*(a_-) - r^*(a_+)) [1 - \sigma(r^*(a_-) - r^*(a_+))] [r(a_+) - r(a_-) - r^*(a_+) + r^*(a_-)]^2$
1030
1031 $\stackrel{(b)}{\geq} \frac{1}{64r_{\max}} \exp(r^*(a_-) - r^*(a_+)) [r(a_+) - r(a_-) - r^*(a_+) + r^*(a_-)]^2.$
1032
1033

1034 Here step (a) follows from Lemma 5; step (b) makes use of the fact that $r^*(a_-) \leq r^*(a_+)$. Hence
1035 we have

1036
$$\begin{aligned} & [r(a_+) - r(a_-) - r^*(a_+) + r^*(a_-)]^2 \\ 1037 & \leq 64r_{\max} \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r(a_+) - r(a_-))). \end{aligned} \quad (\text{C.14})$$

1038
1039

1040 In addition, we have

1041
$$\begin{aligned} \sum_{i=1}^t \text{KL}(\sigma(r^*(a_1^i) - r^*(a_2^i)) \parallel \sigma(r^{(t)}(a_1^i) - r^{(t)}(a_2^i))) & \stackrel{(i)}{\leq} -2[\ell(r^{(t)}, \mathcal{D}^{(t)}) - \ell(r^*, \mathcal{D}^{(t)})] + 2C_3Ar_{\max} \log T \\ 1044 & \stackrel{(ii)}{\leq} 2\alpha_t \gamma_t + 2C_3Ar_{\max} \log T \end{aligned}$$

1045

1046 Here step (i) follows from Lemma 2, while step (ii) utilizes (5.6) and the definition of γ_t (cf. (5.2)).
1047 This immediately implies that

1048
$$\text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t)}(a_+) - r^{(t)}(a_-))) \leq \frac{2\alpha_t \gamma_t + 2C_3Ar_{\max} \log T}{N_t(a_+, a_-)}. \quad (\text{C.15})$$

1049
1050

1051 Therefore for any $1 \leq n_1 < n_2 \leq N_T(a_+, a_-)$, we have

1052
$$\begin{aligned} & \frac{1}{n_2 - n_1} \left(\sum_{n=n_1}^{n_2} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)| \right)^2 \\ 1053 & \stackrel{(i)}{\leq} \sum_{n=n_1}^{n_2} [r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)]^2 \\ 1054 & \stackrel{(ii)}{\leq} 64r_{\max} \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \sum_{n=n_1}^{n_2} \text{KL}(\sigma(r^*(a_+) - r^*(a_-)) \parallel \sigma(r^{(t_n)}(a_+) - r^{(t_n)}(a_-))) \\ 1055 & \stackrel{(iii)}{\leq} \frac{128r_{\max}\alpha_T}{n_1} \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \sum_{n=n_1}^{n_2} \gamma_{t_n} + 128C_3Ar_{\max}^2 \log T \frac{n_2 - n_1}{n_1} \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)}. \end{aligned} \quad (\text{C.16})$$

1056
1057
1058
1059
1060
1061
1062
1063

1064 Here step (i) uses the Cauchy-Schwarz inequality; step (ii) follows from (C.14); step (iii) utilizes
1065 (C.15) and the fact that $\{\alpha_t\}$ is monotonically increasing. Following the same analysis as in (5.3)
1066 and (5.4), we know that

1067
$$\begin{aligned} \sum_{n=n_1}^{n_2} \gamma_{t_n} & \leq \sum_{n=n_1}^{n_2} \xi_{t_n} + \sum_{n=n_1}^{n_2} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)| \\ 1068 & \leq C_1 r_{\max} \sqrt{(n_2 - n_1) \log T} + \sum_{n=n_1}^{n_2} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)|. \end{aligned} \quad (\text{C.17})$$

1069
1070
1071
1072
1073

1074 Taking (C.16) and (C.17) collectively and let $n_2 = 2n_1$, we know that for any $n_1 \leq N_T(a_+, a_-)/2$,

1075
$$\begin{aligned} & \left(\sum_{n=n_1}^{2n_1} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)| \right)^2 \\ 1076 & \stackrel{(iii)}{\leq} 128r_{\max}\alpha_T \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \sum_{n=n_1}^{2n_1} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)| \end{aligned}$$

1077
1078
1079

$$+ 128r_{\max} \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} n_1 \left(\alpha_T C_1 r_{\max} \sqrt{\frac{\log T}{n_1}} + C_3 A r_{\max} \log T \right).$$

This self-bounding relation implies that

$$\begin{aligned}
& \sum_{n=n_1}^{2n_1} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^*(a_+) + r^*(a_-)| \leq 256r_{\max}\alpha_T \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \\
& + \sqrt{256r_{\max} \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} n_1 \left(\alpha_T C_1 r_{\max} \sqrt{\frac{\log T}{n_1}} + C_3 A r_{\max} \log T \right)} \\
& \leq 400r_{\max}\alpha_T \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} + C_1 r_{\max} \sqrt{n_1 \log T} + C_3 n_1 \frac{A r_{\max} \log T}{\alpha_T},
\end{aligned}$$

where the last relation follows from the AM-GM inequality. By using the above relation recursively, we have

$$\begin{aligned}
\zeta(a_+, a_-) &\leq \sum_{k=1}^{\lceil \log T \rceil} \sum_{n=N_T(a_+, a_-)/2^k}^{N_T(a_+, a_-)/2^{k-1}} |r^{(t_n)}(a_+) - r^{(t_n)}(a_-) - r^\star(a_+) + r^\star(a_-)| \\
&\leq C_9 r_{\max} \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T \log T + C_9 r_{\max} \sqrt{N_T(a_+, a_-) \log T} + C_9 N_T(a_+, a_-) \frac{A r_{\max} \log T}{\alpha_T} \\
&\quad \text{(C.18)}
\end{aligned}$$

for some sufficiently large constant $C_9 > 0$. On the other hand, taking (C.18) and (C.13) collectively yields

$$\begin{aligned}
1104 & \\
1105 \quad \zeta(a_+, a_-) \leq C_8 C_9 \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \log T + C_9 r_{\max} \sqrt{N_T(a_+, a_-) \log T} \\
1106 & \\
1107 & \\
1108 & \quad + C_9 N_T(a_+, a_-) \frac{A r_{\max} \log T}{\alpha_T}. \tag{C.19} \\
1109 &
\end{aligned}$$

By putting (C.10), (C.12), (C.18) and (C.19) together, we have

$$\zeta(a_+, a_-) \leq C_6(r_{\max} + \log T) \min \left\{ \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T r_{\max}, \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\} \\ + C_6 \left(\frac{A N_T(a_+, a_-) \log T}{\alpha_T} + \sqrt{T \log T} \right) r_{\max}$$

always holds for some universal constant $C_6 > 0$.

D PROOF OF PROPOSITION 3

Under Assumption 1, we know that for any action pair (a_+, a_-) ,

$$\min \left\{ \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T r_{\max}, \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\} \leq \max \left\{ \tau \alpha_T r_{\max}, (\kappa T)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\}.$$

Therefore we have

$$\begin{aligned} \mathcal{R}(T) &\leq C r_{\max} A^2 \sqrt{T \log T} + C \sum_{t=1}^T \frac{A r_{\max} \log T}{\alpha_t} + 2C(r_{\max} + \log T) A^2 \tau \alpha_T r_{\max} \\ &\quad + C(r_{\max} + \log T) A^2 (\kappa T)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}}. \end{aligned}$$

By taking

$$\alpha_t = A \log T + t^{\frac{1}{\beta+2}} \left(\frac{r_{\max}}{\kappa} \right)^{\frac{\beta}{\beta+2}} \left(\frac{\log T}{A(r_{\max} + \log T)} \right)^{\frac{\beta+1}{\beta+2}},$$

1134 we can achieve

$$\begin{aligned}
 1136 \quad \mathcal{R}(T) &\lesssim (r_{\max} + \log T) A^3 \tau r_{\max} \log T + r_{\max} A^2 \sqrt{T \log T} \\
 1137 &\quad + (r_{\max} + \log T)^{\frac{\beta+1}{\beta+2}} r_{\max}^{\frac{2}{\beta+2}} \kappa^{\frac{\beta}{\beta+2}} A^{\frac{2\beta+3}{\beta+2}} T^{\frac{\beta+1}{\beta+2}} (\log T)^{\frac{1}{\beta+2}} \\
 1138 &\quad + (r_{\max} + \log T)^{\frac{1}{\beta+2}} A^{\frac{\beta+3}{\beta+2}} \tau r_{\max}^{\frac{2\beta+2}{\beta+2}} \kappa^{-\frac{\beta}{\beta+2}} (\log T)^{\frac{\beta+1}{\beta+2}} T^{\frac{1}{\beta+2}} \\
 1139 &\quad + (r_{\max} + \log T) A^{\frac{2\beta+3}{\beta+1}} \kappa^{\frac{\beta}{\beta+1}} (\log T)^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} T^{\frac{\beta}{\beta+1}} \\
 1140 &\lesssim \tau A^3 r_{\max}^2 \log^2 T + T^{\frac{\beta+1}{\beta+2}} \kappa^\beta r_{\max}^2 A^3 \tau \log^2 T.
 \end{aligned}$$

1144 E ANOTHER ASSUMPTION AND THE REGRET BOUND

1145 As an alternative to Assumption 1, we can also impose the following assumption to capture the
1146 relation between human preference π_{HF} and the reference policy π_{ref} .

1147 **Assumption 2.** *There exists some quantity $\mu > 0$ such that, for any action pair (a_+, a_-) ,*

$$\begin{aligned}
 1150 \quad \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} &\leq \mu \frac{\pi_{\text{ref}}(a_+)}{\pi_{\text{ref}}(a_-)}.
 \end{aligned}$$

1153 The quantity μ measures the deviation of human preference from the reference policy. Under
1154 Assumption 2, we have

$$\begin{aligned}
 1156 \quad \min \left\{ \frac{\pi_{\text{HF}}(a_+)}{\pi_{\text{HF}}(a_-)} \alpha_T r_{\max}, \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\} \\
 1157 &\leq \min \left\{ \mu \frac{\pi_{\text{ref}}(a_+)}{\pi_{\text{ref}}(a_-)} \alpha_T r_{\max}, \left(T \frac{\pi_{\text{ref}}(a_-)}{\pi_{\text{ref}}(a_+)} \right)^{\frac{\beta}{\beta+1}} \alpha_T^{\frac{1}{\beta+1}} r_{\max}^{\frac{1}{\beta+1}} \right\}. \\
 1158 &\leq (\mu T)^{\frac{\beta}{2\beta+1}} (\alpha_T r_{\max})^{\frac{\beta+1}{2\beta+1}}.
 \end{aligned}$$

1163 Putting the above relation with (4.1), we have

$$\begin{aligned}
 1165 \quad \mathcal{R}(T) &\lesssim r_{\max} A^2 \sqrt{T \log T} + \sum_{t=1}^T \frac{A r_{\max} \log T}{\alpha_t} + A^2 \alpha_T r_{\max}^2 \\
 1166 &\quad + (r_{\max} + \log T) A^2 (\mu T)^{\frac{\beta}{2\beta+1}} (\alpha_T r_{\max})^{\frac{\beta+1}{2\beta+1}}.
 \end{aligned}$$

1169 By taking

$$\alpha_t = A + t^{\frac{\beta+1}{3\beta+2}} \left(\frac{r_{\max}}{\mu} \right)^{\frac{\beta}{3\beta+2}} \left(\frac{\log T}{A(r_{\max} + \log T)} \right)^{\frac{2\beta+1}{3\beta+2}},$$

1172 we have

$$\mathcal{R}(T) \lesssim T^{\frac{2\beta+1}{3\beta+2}} \mu^{\frac{\beta}{3\beta+2}} \text{poly}(A, r_{\max}, \log T).$$

1176 F TECHNICAL LEMMAS

1178 **Lemma 5.** *For any $x, \delta \in \mathbb{R}$, we have*

$$\text{KL}(\sigma(x) \| \sigma(x + \delta)) \geq \frac{1}{4} \sigma(x) (1 - \sigma(x)) \min\{|\delta|, \delta^2\}.$$

1183 *Proof.* Let $f_x(t) := \text{KL}(\sigma(x) \| \sigma(x + t))$. We have

$$\begin{aligned}
 1184 \quad f_x(t) &= \sigma(x) \log \frac{\sigma(x)}{\sigma(x + t)} + (1 - \sigma(x)) \log \frac{1 - \sigma(x)}{1 - \sigma(x + t)} \\
 1185 &= \sigma(x) \log \left(\frac{\sigma(x)}{1 - \sigma(x)} \cdot \frac{1 - \sigma(x + t)}{\sigma(x + t)} \right) + \log \frac{1 - \sigma(x)}{1 - \sigma(x + t)}
 \end{aligned}$$

$$= \log \frac{1 + \exp(x + t)}{1 + \exp(x)} - \sigma(x)t = \log(1 + \sigma(x)(e^t - 1)) - \sigma(x)t.$$

Then we have

$$f'_x(t) = \frac{\sigma(x)e^t}{1 + \sigma(x)(e^t - 1)} - \sigma(x) = \frac{\sigma(x)(1 - \sigma(x))(e^t - 1)}{1 + \sigma(x)(e^t - 1)}.$$

For any $t > 0$, we can check that

$$f'_x(t) > \sigma(x)(1 - \sigma(x))(1 - e^{-t}) \geq \frac{1}{2}\sigma(x)(1 - \sigma(x))\min\{t, 1\},$$

and for any $t \in (0, 1)$ we have

$$f'_x(t) < \sigma(x)(1 - \sigma(x))(e^t - 1) \leq 2\sigma(x)(1 - \sigma(x))t.$$

This immediately implies that for $\delta > 0$,

$$\begin{aligned} \text{KL}(\sigma(x)\|\sigma(x + \delta)) &= f_x(\delta) - f_x(0) = \int_0^\delta f'_x(t)dt \\ &\geq \frac{1}{2}\sigma(x)(1 - \sigma(x)) \int_0^\delta \min\{t, 1\}dt \\ &\stackrel{(a)}{\geq} \frac{1}{4}\sigma(x)(1 - \sigma(x))\min\{\delta, \delta^2\}. \end{aligned}$$

Here step (a) holds since $\int_0^\delta \min\{t, 1\}dt = \delta^2/2$ for $\delta \leq 1$, and $\int_0^\delta \min\{t, 1\}dt = \delta - 1/2 \geq \delta/2$ for $\delta > 1$.

For $\delta < 0$, we can use the same argument to show that

$$\text{KL}(\sigma(x)\|\sigma(x + \delta)) \geq \frac{1}{4}\sigma(x)(1 - \sigma(x))\min\{-\delta, \delta^2\}.$$

This completes the proof. \square

The following lemma provides a user-friendly version of Freedman's inequality (the Bernstein inequality for martingale differences) (Freedman, 1975; Tropp, 2011).

Lemma 6. Consider a filtration $\{\mathcal{F}_i\}_{i \geq 0}$ and random variables $\{X_i\}_{i \geq 1}$ obeying

$$|X_i| \leq R \quad \text{and} \quad \mathbb{E}[X_i|\mathcal{F}_{i-1}] = 0 \quad \text{for all } i \geq 1.$$

Define $W_n = \sum_{i=1}^n \mathbb{E}[X_i^2|\mathcal{F}_{i-1}]$, and suppose that $W_n \leq \sigma^2$ holds deterministically for some given quantity $\sigma > 0$. Then for any positive integer $m \geq 1$, with probability exceeding $1 - \delta$ we have

$$\left| \sum_{i=1}^n X_i \right| \leq \sqrt{8 \max \left\{ W_n, \frac{\sigma^2}{2^m} \right\} \log \frac{2m}{\delta}} + \frac{4}{3}R \log \frac{2m}{\delta}.$$

Proof. See Li et al. (2021, Section A). \square

1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241