General Exploratory Bonus for Optimistic Exploration in RLHF

Anonymous Author(s)
Affiliation
Address
email

Abstract

Optimistic exploration is central to improving sample efficiency in reinforcement learning with human feedback, yet existing exploratory bonus methods to incentivize exploration often fail to realize optimism. We provide a theoretical analysis showing that current formulations, under KL or α -divergence regularization, unintentionally bias exploration toward high-probability regions of the reference model, thereby reinforcing conservative behavior instead of promoting discovery of uncertain regions. To address this pitfall, we introduce the **General Exploratory Bonus** (**GEB**), a novel theoretical framework that provably satisfies the optimism principle. GEB counteracts divergence-induced bias via reference-dependent reward regulation and unifies prior heuristic bonuses as special cases, while extending naturally across the full α -divergence family. Empirically, GEB consistently outperforms baselines on alignment tasks across multiple divergence settings and large language model backbones. These results demonstrate that GEB offers both a principled and practical solution for optimistic exploration in RLHF. Code is available here.

1 Introduction

Despite the acknowledged significance of online exploration for reinforcement learning with human feedback (RLHF) [1, 2, 3], there remains a paucity of theoretical frameworks governing *how to explore*. As shown in Fig. 1 (1, top), standard online RLHF algorithms [4, 5, 6] generally rely on passive exploration, *i.e.*, the stochasticity of the policy itself to generate responses, with no mechanism to incentivize novelty or diversity. As a result, this approach can be notoriously sample-inefficient. When the optimal behavior resides in low-probability regions, passive exploration is unlikely to discover it, leading to policies that remain trapped around local optima.

To address this, some works [7, 8, 9, 10, 11] have attempted to devise sample-efficient algorithms, inspired by the principle *optimism in the face of uncertainty*. As illustrated in Fig. 1 (2, top), the principle aims to generate responses with high epistemic uncertainty, thus encouraging data collection in unexplored regions for further training. To operationalize the principle, recent research [12, 13, 14] encourages exploration by adding *exploratory bonuses* to the reward modeling, which is practically optimizable for large language models. These methods intend to artificially inflate rewards in underexplored regions, nudging the policy toward more informative data collection.

Unfortunately, our theoretical analysis in Section 3 reveals a fundamental pitfall: under the common KL-regularized RLHF, the existing theoretical framework of exploratory bonuses fails to satisfy optimism. In particular, we prove that existing bonus formulations can undesirably drive the policy π toward the reference policy π_{ref} due to the divergence regulation in the exploratory bonus, and the induced bonus actually biases exploration toward high-probability regions of the reference model. As illustrated in Fig. 1 (II, bottom), the bonus disproportionately amplifies rewards for regions already well-covered by π_{ref} , thereby reinforcing conservative behavior rather than driving exploration into

uncertain regions. This failure is not confined to KL-divergence; we further extend our analysis to the more general α -divergence family and prove that the same collapse persists across a wide range of divergence-regularized objectives. Thus, while existing approaches appear to encourage exploration, they in fact undermine the very principle of optimism they aim to realize.

Motivated by these failures, we propose a new framework, **General Exploratory Bonus** (GEB), 41 which theoretically unifies existing approaches while provably satisfying optimism (Section 4). 42 GEB corrects the failure modes of prior approaches by directly introducing a reference-dependent 43 regulation into the reward. This adjustment offsets the undesired conservatism induced by divergence 44 regularization, allowing the exploratory bonus to satisfy optimism—it increases the probability of 45 responses rarely sampled to pursue potentially more preferred answers, as shown in Fig. 1 (III, 46 bottom). Importantly, GEB provides a unified formulation: prior heuristic exploratory bonuses can be 47 reinterpreted as special cases, and the framework naturally extends to the full class of α -divergences. 48 Beyond correcting the theoretical shortcomings, GEB remains practically implementable—it can be 49 seamlessly integrated into the standard iterative RLHF loop without additional sampling cost.

We validate GEB on a large-scale alignment task across different divergences and model backbones. Empirically, GEB consistently yields stronger alignment compared to its counterpart of passive exploration. For example, the three GEB variants that we consider generally outperform the iterative f-DPO [15] across different divergence regulations, while the most performant variant surpasses several existing optimistic exploration methods that incorporate exploratory bonuses [12, 13, 14]. By analyzing the distribution of sampled responses, we validate that GEB can successfully encourage sampling in the region of small $\pi_{\rm ref}$, thereby effectively achieving optimistic exploration.

We summarize our main contributions:

- 1. We formally prove that the existing theoretical framework of exploratory bonuses under KL and α -divergence regularization fails to achieve optimistic exploration.
- We introduce General Exploratory Bonus (GEB), a novel theoretical framework of optimistic exploration for RLHF that provably satisfies the optimism principle and unifies prior heuristic bonuses.
- 3. We empirically validate GEB on LLM alignment tasks, showing improved performance and broad applicability across multiple divergence families.

2 Preliminaries

58

59

60

61

62 63

64

65

Iterative online RLHF. The effectiveness of iterative online RLHF [16, 17] has been validated in 67 many real-world systems such as Claude [18] and LLaMA-series [19, 20]. The algorithm proceeds 68 in T rounds, with each round having two steps: (1) the π_t is learned with the current dataset 70 \mathcal{D}_t , and then samples $x \sim \rho$, $(y_1, y_2) \sim \pi_t(\cdot | x)$; (2) Human evaluator annotate the preference of $(x, y_1, y_2) \to (x, y^w, y^l)$ to form \mathcal{D}_{t+1} , where the prompt x is sampled from an independent 71 distribution ρ , response y_1, y_2 are two response sampled from the policy of the t-th iteration π_t . 72 When computing π_t with dataset D_t in the step (1) of each iteration, a reward function $r_t(x,y)$ 73 is first learned from a collected human preference data $\mathcal{D}_t = \{(x, y^w, y^l)\}$, where y^w, y^l denote 74 75 the preferred and dispreferred response to x, respectively. Reward modeling typically follows the Bradley-Terry objective [21]: 76

$$r_t = \arg\min_r \mathcal{L}_{BT}(\mathcal{D}_t, r) = \arg\min_r \mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}_t} - \log[\sigma(r(x, y^w) - r(x, y^l))], \tag{1}$$

where σ denotes the sigmoid function. Given the learned reward function r_t , the policy π_t is then updated to maximize the expected reward, often with a KL-regularization as follows

$$\pi_t = \arg\max_{\pi} \mathcal{J}_{\beta, \text{KL}}(\pi, r_t) = \arg\max_{\pi} \mathbb{E}_{x \sim \rho, y \sim \pi(\cdot | x)} r_t(x, y) - \beta \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}}), \tag{2}$$

where $\beta > 0$ is a hyperparameter, and $\pi_{\rm ref}$ is the reference model.

Sample inefficiency of iterative online RLHF. Online sampling for standard online RLHF algorithms is carried out passively, relying solely on the inherent randomness of the LLM policy. However, if the policy places a small probability mass on the optimal action, passive exploration may fail to ever explore this action. Theoretical analyses [22, 17] and empirical evidence [23, 14] present that the

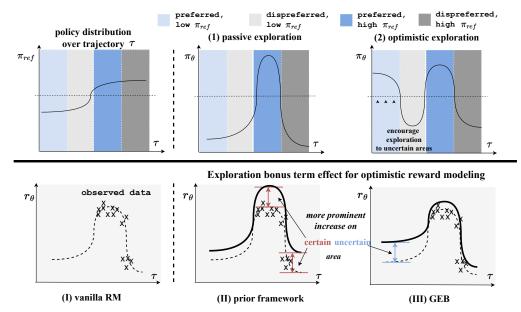


Figure 1: The upper part compares passive exploration and optimistic exploration. Optimistic exploration stimulates the trajectories τ of small $\pi_{\rm ref}$ (seldom visited/uncertain). While passive exploration sticks to the high- π_{ref} region, failing to approach global optima. The lower part contrasts the effect of the exploration bonus term in optimistic reward modeling between prior works and our GEB. Prior works often emphasize rewards in frequently visited regions, which constrains exploration within certain areas. In contrast, our GEB amplifies rewards in seldom-visited regions, thereby encouraging further sampling in uncertain areas and successfully achieving optimistic exploration.

passive approach fails to sufficiently explore the prompt-response space. Particularly, Xie et al. [13] demonstrate that the sample complexity can be exponential in $1/\beta$ for passive exploration, which 85 is unacceptable in the small- β regime. Therefore, several works [13, 14, 12] propose exploratory 86 bonuses to implement optimistic exploration for efficient sampling. However, in the next section, we 87 will show that prior formulations cannot provably achieve optimism. 88

Exploratory Bonus and How It Can Fail 89

In this section, we will first provide the iterative online RLHF formulation with an exploratory bonus 90 (Section 3.1). We then theoretically prove that the existing formulation can fail to achieve optimistic 91 exploration under both KL-constrained RLHF (Section 3.2) and a more general α -divergence-92 regularized RLHF (Section 3.3), motivating our proposed method in Section 4. 93

Exploratory Bonus 94

To improve the sample efficiency of iterative online RLHF, recent works [12, 14] introduce exploratory 95 bonuses, which try to encourage optimistic exploration. These approaches modify the standard loop 96 by adding an exploratory bonus term $\mathcal{L}_{\text{bonus}}$. Specifically, in the t-th iteration, the reward model r_t 97 and policy π_t are optimized by 98

$$r_{t} = \arg\min_{r} \left[\mathcal{L}_{BT}(\mathcal{D}_{t}, r) - \kappa \mathcal{L}_{\text{bonus}}(r) \right],$$

$$\pi_{t} = \arg\max_{\pi} \mathcal{J}_{\beta, \text{KL}}(\pi, r_{t}) = \arg\max_{\pi} \mathbb{E}_{x \sim \rho, y \sim \pi(\cdot | x)} r_{t}(x, y) - \beta \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}}),$$
(4)

$$\pi_t = \arg\max_{\pi} \mathcal{J}_{\beta, \mathsf{KL}}(\pi, r_t) = \arg\max_{\pi} \mathbb{E}_{x \sim \rho, y \sim \pi(\cdot \mid x)} r_t(x, y) - \beta \mathbb{D}_{\mathsf{KL}}(\pi \mid \pi_{\mathsf{ref}}), \tag{4}$$

where $\kappa > 0$ is a hyperparameter. By Eq. 3, the reward model r_t should not only fit the observed data 99 in \mathcal{D}_t , but also learn to maximize the bonus term $\mathcal{L}_{\text{bonus}}(r)$. 100

To achieve optimistic exploration, the bonus term is expected to stimulate probability increase more prominently in unexplored areas. Formally, we have the following definition:

Definition 3.1 (Optimism condition for exploration bonus) When a reward model r and a policy π is computed with Eq. 3 and Eq. 4, the exploratory bonus $\mathcal{L}_{bonus}(r)$ achieves optimism, if

$$\frac{\partial^2 \mathcal{L}_{bonus}(r)}{\partial \pi(y|x)\partial \pi_s(y|x)} < 0, \tag{5}$$

where $\pi_s(y|x)$ is an ideal sampling distribution for response at the current iteration.

To interpret the optimism condition of the exploration bonus term in Definition. 3.1, we consider the policy-reparameterized reward model $r(x,y):=r(\pi)$, which can be derived from the closed-form solution of Eq. 2 as $\pi(y|x)=\frac{\exp(r(x,y)/\beta)}{Z(x)}$, where $Z(x)=\mathbb{E}_{y\sim\pi_{\rm ref}}\exp(r(x,y)/\beta)$ is a normalization function. This yields the reward model expressed via the policy [24]:

$$r(\pi) = \beta \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} + \log Z(x).$$

Implication. Eq. 5 requires the gradient $\partial \mathcal{L}_{\text{bonus}}(r(\pi))/\partial \pi(y|x)$ to be negatively correlated with the sampling probability π_s . In other words, responses rarely sampled under π_s (i.e., uncertain or underexplored outputs) should receive a larger ascending of the policy distribution π , i.e., larger $\partial \mathcal{L}_{\text{bonus}}(r(\pi))/\partial \pi(y|x)$. In practice, π_s can be substituted by the reference model or intermediate checkpoints. We adopt the commonly used π_{ref} as π_s in our following demonstration.

111 3.2 Failure Under KL-constrained RLHF

Previous works, including Zhang et al. [12] and Cen et al. [14], formulate the exploratory bonus with $\mathcal{L}_{\mathrm{bonus}}(r) = \max_{\pi} \mathcal{J}_{\beta,KL}(\pi,r)$. In this case, optimizing exploratory bonus in Eq. 3 becomes a min-max bi-level objective as $\min_r -\kappa \max_{\pi} [\mathbb{E}_{x,y \sim \pi} r(x,y) - \beta \mathbb{D}_{\mathrm{KL}}(\pi \| \pi_{\mathrm{ref}})]$. Intuitively, they intend to make r not only fit the observed data by \mathcal{L}_{BT} , but also have a larger reward in unobserved regions by maximizing the $\max_{\pi} \mathbb{E}_{x,y \sim \pi} r(x,y)$ in $\mathcal{L}_{\mathrm{bonus}}(r)$. Here, we theoretically show that such formulations can suffer from optimism failures under KL-regularized RLHF.

Lemma 3.1 (Optimism failure under KL-divergence.) Let $r_1 = \arg\min_r \mathcal{L}_{BT}(\mathcal{D}, r)$ be a reward model trained with the vanilla BT loss, and let $r_2 = \arg\min_r [\mathcal{L}_{BT}(\mathcal{D}, r) - \kappa \max_{\pi} \mathcal{J}_{\beta, KL}(\pi, r)]$ be a reward model trained with an additional exploratory bonus. If the policy is optimized via Eq. 4, then r_1 and r_2 yield the same set of policies.

122 Proof First, the inner maximization of the bonus term admits a closed-form solution, $\pi^*(y|x)=123$ $\pi_{ref}(y|x)e^{\frac{r(x,y)}{\beta}}/Z(x)$ where $Z(x)=\mathbb{E}_{y\sim\pi_{ref}(\cdot|x)}e^{\frac{r(x,y)}{\beta}}$ is a normalization term. Substituting this solution reduces the bi-level training objective of r_2 to a single-level form:

$$r_2 = \arg\min_{r} \left[\mathcal{L}_{BT}(\mathcal{D}, r) - \kappa \mathbb{E}_{x \sim \rho} \beta \log \mathbb{E}_{y \sim \pi_{\text{ref}}} e^{\frac{r(x, y)}{\beta}} \right]. \tag{6}$$

As shown in Rafailov et al. [24], the log-ratio $\beta \log \pi_{\theta}(y|x) - \beta \log \pi_{\text{ref}}(y|x)$ represents the same class of the original reward function r through Eq. 4, thus the reward r_1, r_2 can be reparameterized by the log-ratio. Plugging this into Eq. 6 yields

$$\arg\min_{\pi} \mathcal{L}_{dpo}(\mathcal{D}, \pi) - \kappa \mathbb{E}_{x \sim \rho} \beta \log \mathbb{E}_{y \sim \pi_{\text{ref}}(\cdot|x)} \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)}. \tag{7}$$

Since the second term equals 0, the reparameterized Eq. 6 is exactly the vanilla DPO loss, which is the same as the reparameterized training objective of r_1 . Thus, the exploratory bonus in the reward training objective has no effect on the final policy set.

The lemma proves that incorporating the exploratory bonus $\mathcal{L}_{\text{bonus}}(r) = \max_{\pi} \mathcal{J}_{\beta,\text{KL}}(\pi,r)$ into the reward training objective *fails to induce the policy model to sample from low-\pi_{\textit{ref}}(y|x) regions, i.e., unexplored responses.* In other words, the bonus term is ineffective at inducing optimism. We next extend the result beyond KL regularization to a more general class of α -divergence family.

3.3 Generalization to α -divergence-constrained RLHF

135

In this subsection, we theoretically show that the failure of optimism can broadly extend to the α -divergence class. Many common divergences, such as reverse KL-divergence, Hellinger distance, and forward KL-divergence, are special cases of α -divergence.

Definition 3.2 (α -divergence class) An α -divergence is a certain type of function D(p|q) = $\int f(\frac{dp}{dq})dq$ that measures the difference between two probability distributions p and q, where

$$f(x) = \frac{x^{\alpha} - \alpha x - (1 - \alpha)}{\alpha (1 - \alpha)},$$

and α is a hyperparameter typically with $0 \le \alpha \le$

Lemma 3.2 (Optimism failure under α -divergence.) Consider objective $\mathcal{J}_{\beta,f}(\pi,r)$ 140 $\mathbb{E}_{x \sim \rho, y \sim \pi(y|x)} r(x, y) + \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{ref}(y|x)} f(\frac{\pi(y|x)}{\pi_{ref}(y|x)})$, where f belongs to α -divergence class. If a reward is trained with $r = \arg\min_{r} [\mathcal{L}_{BT}(\mathcal{D}, r) - \kappa \mathcal{L}_{bonus}]$ and a policy π is updated by $\arg\max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$ with $\mathcal{L}_{bonus} = \max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$, the gradient of the bonus satisfies 141 $rac{\partial^2 \mathcal{L}_{bonus}(r(\pi))}{\partial \pi \partial \pi_{ref}} \geq 0$, which means \mathcal{L}_{bonus} encourage trajectories with large π_{ref} more strongly, in 144 contradiction to the optimism principle (Definition 3.1). 145

Proof For a RL objective $\mathcal{J}_{\beta,f}(\pi,r)$, the relation between the optimal policy π_f^* and the reward r146 can be formulated as follows, 147

$$\pi_f^* = \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) (f')^{-1} (r/\beta), \quad r(x,y) = \beta f'(\frac{\pi^*(y|x)}{\pi_{\text{ref}}(y|x)} Z(x)), \tag{8}$$

where Z(x) is a normalization term and $(f')^{-1}$ is the inverse function of f'. The bi-level objective can 148 be similarly transformed to a single level one by canceling the inner maximization \max_{π} by Eq. 8. The 149 single-level objective can be written as $r_t = \arg\min_r \mathcal{L}_{BT}(\mathcal{D}, r) - \kappa \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \frac{1}{Z(x)} (f')^{-1} (\frac{r(x,y)}{\beta}) \cdot r(x,y) - \beta f(\frac{1}{Z(x)}(f')^{-1} (\frac{r(x,y)}{\beta}))$. Since the policy is computed by $\arg\max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$, the reward 150 151 can be reparameterized by the policy with Eq. 8, which fortunately cancels Z(x). Then the optimistic reward-modeling objective can be reparameterized as 152 153

$$\arg\min_{\pi} \mathcal{L}_{dpo}(\mathcal{D}, \pi) - \kappa \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \left[\frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} f'(\frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)}) - \beta f(\frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)}) \right]. \tag{9}$$

Since for α -divergence, $f(u) = \frac{u^{\alpha} - \alpha u - (1 - \alpha)}{\alpha(\alpha - 1)}$, the partial derivative of Eq. 9 is $(\frac{\pi_{\text{ref}}}{\pi})^{1 - \alpha}$, which 154 induces positively correlated gradients w.r.t. π and π_{ref} when $0 \le \alpha < 1$, and is a constant when 155 $\alpha = 1$, hence contradictory to the optimism defined in Definition 3.1. 156

According to Lemma. 3.2, we enumerate several ex- Table 1: Realized exploratory bonus under difploratory bonus under different α -divergence in Ta- ferent divergence classes when $\mathcal{L}_{\text{bonus}}(r)$ ble 1. The listed bonuses are simplified by removing constant coefficients and bias. The listed exploratory bonuses generally force the policy model to maximize the possibility of trajectories sampled by the reference model, not the underexplored ones. We further prove that it actually drives π to collapse toward $\pi_{\rm ref}$ and that the failure extends beyond α -divergence to other f-divergences.

157

158

159

160

161

162

163

164

165

166

167

168

169 170

 $\max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$.

f	exploratory bonus
reverse KL	constant
forward KL	$\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}(y x)} \log \frac{\pi(y x)}{\pi_{\text{ref}}(y x)}$
Hellinger distance	$\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}(y x)} \sqrt{\frac{\pi(y x)}{\pi_{\text{ref}}(y x)}}$

Theorem 3.3 (Optimism failure beyond α **-divergence.)** When f belongs to f-divergence, and the reward function is obtained by $\hat{r} = \arg\min_r [\mathcal{L}_{BT}(\mathcal{D}_t, r) - \kappa \max_{\pi} \mathcal{J}_{\beta, f}(\pi, r)]$ and the policy is updated by $\arg \max_{\pi} \mathcal{J}_{\beta,f}(\pi, r_t)$, the bonus term $-\kappa \max_{\pi} \mathcal{J}_{\beta,f}(\pi, r)$ induces the policy model π to coincide with π_{ref} when xf''(x) is a monotone function.

The detailed proofs are in Appendix B.1. The monotone increase of xf''(x) can be satisfied by a 171 broader divergence class besides α -divergence, including JS-divergence and Pearson χ^2 . 172

Intuitive understanding. Optimization of exploratory bonus in Eq. 3 is a min-max bi-level objec-173 tive as $\min_r -\kappa \max_{\pi} [\mathbb{E}_{x,y \sim \pi} r(x,y) - \beta \mathbb{D}_{KL}(\pi || \pi_{ref})]$. Due to the inner maximization \max_{π} , the 174 divergence constraint implicitly makes π close to π_{ref} to avoid the divergence penalty. Considering 175 the outer minimization \min_r , r is forced to provide large rewards on region of high π to maximize 176 the reward expectation. Their combination implicitly makes r focus more on region of large π_{ref} . 177 Since samples with large π_{ref} is easily rolled out from scratch, previous exploratory bonuses merely 178 concentrate sampling on regions that are already easy to explore, contradictory to the optimism principle, which requires encouraging responses y rarely sampled by the reference model.

4 General Exploratory Bonus with Optimism Principle

181

Motivated by the failure of the existing exploratory bonus, we now propose a novel framework, *General Exploratory Bonus* (**GEB**), and prove that it achieves optimism. We further show that prior exploratory bonuses—and broader variants—emerge as special cases of our formulation.

Formulation of a novel exploratory bonus. As shown in Section 3, the failure of existing bonuses 185 arises because the divergence constraints in $\max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$ force π to align with π_{ref} , biasing 186 exploration toward high- π_{ref} regions. To achieve optimistic exploration, the optimal π must instead 188 counteract this regularization and move away from π_{ref} . Our key idea is to introduce an additional reference-dependent regulation into the reward, which offsets the influence of divergence regular-189 ization. The resulting exploratory bonus takes the form $-\kappa \max_{\pi} J_{\beta,f}(\pi, R(r, \pi_{\text{ref}}))$. Note that the 190 formulation of $R(\cdot, \cdot)$ can be diverse. Since the reward $R(r, \pi_{\text{ref}})$ is now explicitly dependent on 191 π_{ref} , the inner optimal policy of the bi-level problem is $\frac{1}{Z_R(x)}\pi_{\text{ref}}(f')^{-1}(\frac{R(r,\pi_{\text{ref}})}{\beta})$, where $Z_R(x)$ is a 192 normalization term. Unlike previous cases, the optimal policy is no longer guaranteed to be positively 193 correlated with $\pi_{\rm ref}$, enabling the policy to deviate from the reference distribution. 194

As in Lemma 3.2, we can substitute the inner π by the closed-form solution, and then utilize the reward reparameterization of α -divergence [25] as $r=f'(\pi/\pi_{\rm ref})$ for a divergence instance f to obtain the reparameterized exploratory bonus as

$$\mathcal{L}_{\text{bonus}}(r(\pi)) = \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}(\cdot|x)} \left[\frac{u}{Z_R(x)} f'(u) - f(\frac{u}{Z_R(x)}) \right], \tag{10}$$

where $u=(f')^{-1}(R(\pi_{\mathrm{ref}}(y|x),\frac{\pi(y|x)}{\pi_{\mathrm{ref}}(y|x)})/\beta)$ and $Z_R(x)=\mathbb{E}_{y\sim\pi_{\mathrm{ref}}(\cdot|x)}u$. Since the domain of divergence class is generally $(0,+\infty)$ while there are no additional constraints on the formulation of R,u can be flexibly formulated with π and π_{ref} unless $u\geq 0$.

Equivalence to a practical objective. In our proposed exploratory bonus, the normalization term $Z_R(x)$ in Eq. 10 cannot be canceled. Fortunately, we prove the following lemma in Appendix E to show the equivalence between the two training objectives, one with and the other without $Z_R(x)$, which helps transform the objective to a succinct formulation for analyses and practical use.

Lemma 4.1 Denote two objectives as $h(u) = \mathbb{E}_{x \sim \rho, y \sim \pi_{ref}} u f'(u) - f(u)$ and $\hat{h}(u) = \mathbb{E}_{x \sim \rho, y \sim \pi_{ref}} \frac{u}{Z(x)} f'(u) - f(\frac{u}{Z(x)})$ where u is a function with π and π_{ref} . When $[f'(u) + u f''(u) - f'(\frac{u}{Z(x)})]/[Z(x)uf''(u)] = \Lambda(x)$ is constant in y and $\Lambda(x) > 0$, minimizing the two objectives $\min_{\pi} -h(u)$ and $\min_{\pi} -\hat{h}(u)$ induce the same class of policies.

GEB successfully achieves optimism. Building on Lemma 4.1, we now prove that our proposed framework indeed achieves the optimism requirement.

Theorem 4.2 For each iteration of online RLHF, if the policy is updated by $\pi = \arg\max_{\pi} J_{\beta,f}(\pi,r)$ while its reward is trained with $r = \arg\min_{r}[\mathcal{L}_{BT}(\mathcal{D},r) - \kappa\mathcal{L}_{bonus}]$. When f belongs to α -divergence class, and $\mathcal{L}_{bonus} = \max_{\pi} J_{\beta,f}(\pi,R(r,\pi_{ref}))$, denote $u(\pi,\pi_{ref}) = (f')^{-1}\left(R\left((f')^{-1}(\frac{\pi}{\pi_{ref}}),\pi_{ref}\right)/\beta\right)$, Then, $\frac{\partial^2 \mathcal{L}_{bonus}}{\partial \pi \partial \pi_{ref}} \leq 0$ if $\forall (x,y)$; $\frac{\partial u}{\partial \pi} + \pi_{ref} \frac{\partial^2 u}{\partial \pi \partial \pi_{ref}} + \frac{(\alpha-1)\pi_{ref}}{u} \frac{\partial u}{\partial \pi} \frac{\partial u}{\partial \pi_{ref}} < 0$ and $u > \alpha$, where α is a hyperparameter that defines α -divergence (Definition 3.2).

Proof First, substituting the optimal solution of the inner maximization and utilizing reward
 reparameterization, we obtain the training objective as in Eq. 10. By Lemma 4.1, this can be
 equivalently expressed as

$$\mathcal{L}_{\text{bonus}} = \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \Big[u f'(u) - f(u) \Big]. \tag{11}$$

For α -divergences, the conditions in Lemma 4.1 are satisfied. Since $f'(u) + uf''(u) - f'(\frac{u}{Z}) = u^{\alpha-1}(Z^{1-\alpha} - \alpha)/(1-\alpha)$ and $uf''(u) = u^{\alpha-1}$, the fraction $\Lambda(x) = (Z^{1-\alpha} - \alpha)/Z(1-\alpha)$ in Lemma 4.1 is independent of y. Since $u > \alpha$, $Z_R = \mathbb{E}_{y \sim \pi_{\rm ref}(\cdot|x)} u > 0$, thus $\Lambda(x) > 0$ is also satisfied. Finally, the mixed second-order derivative of Eq. 11 is computed as

$$\frac{\partial^{2} \mathcal{L}_{\text{bonus}}}{\partial \pi \partial \pi_{\text{ref}}} = \beta \mathbb{E}_{x \sim \rho} \sum_{u} u^{\alpha - 1} \left(\frac{\partial u}{\partial \pi} + \pi_{\text{ref}} \frac{\partial^{2} u}{\partial \pi \partial \pi_{\text{ref}}} + \frac{(\alpha - 1)\pi_{\text{ref}}}{u} \frac{\partial u}{\partial \pi} \frac{\partial u}{\partial \pi_{\text{ref}}} \right) < 0, \tag{12}$$

Table 2: General exploratory bonus under different divergence classes and design of u. Note all $u>\alpha$ when $0<\pi<1$. The presented bonus is simplified by removing constant coefficients or biases.

$\mathcal{L}_{ ext{bonus}}$ u	$1 + \alpha - \pi$	$1/\pi$	$\operatorname{arctanh}(1-\pi) + \alpha$
reverse KL forward KL Hellinger Distance	$ \begin{vmatrix} \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} - \pi(y x) \\ \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \log(1 - \pi(y x)) \\ \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \sqrt{1.5 - \pi(y x)} \end{vmatrix} $	$\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \frac{1}{\pi(y x)}$ $\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} - \log \pi(y x)$ $\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \frac{1}{\sqrt{\pi(y x)}}$	$ \begin{vmatrix} \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \operatorname{arctanh}(1 - \pi(y x)) \\ \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \log \operatorname{arctanh}(1 - \pi(y x)) \\ \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} \sqrt{\operatorname{arctanh}(1 - \pi(y x)) + 0.5} \end{vmatrix} $

which achieves the optimism defined in Definition 3.1.

In our formulation, u can be flexibly defined in terms of π , π_{ref} as long as it satisfies the derivative condition in Theorem 4.2 and $u > \alpha$. This flexibility highlights the extensibility of our framework. In particular, when u depends only on π , any function with $u > \alpha$ and negative correlation with π qualifies. In Table 2, we list several such choices of u, along with their corresponding reparameterized exploratory bonus under three different α -divergences.

From a practical standpoint, since \mathcal{L}_{bonus} is expressed as an expectation over π_{ref} , it does not require additional sampling and can be seamlessly integrated into iterative online RLHF. To avoid unintended decreases in the likelihood of preferred responses, however, we follow Chen et al. [23] and restrict computation of the bonus on rejected responses to ensure that the probability of preferred responses continues to increase.

Prior exploratory bonuses are encompassed within GEB. Although we have shown that existing theoretical formulations of $\mathcal{L}_{\text{bonus}}$ fail to guarantee optimism, many practical implementations have nevertheless been effective through various approximations and adaptations. These approximations and adaptations are generally inextensible beyond the reverse KL divergence (detailed in Appendix B.2). In this subsection, we show that these practical implementations can be naturally subsumed into our GEB framework, and even broader objectives can be reinterpreted as instances of optimistic exploration. For example, Zhang et al. [12] and Xie et al. [13] finally implement their exploratory bonus as $\kappa \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}}(y|x) \log \pi(y|x)$, which belongs to GEB when $u(\pi) = -\log \pi + 1$ and f is KL-divergence. Similarly, Cen et al. [14] implement the exploratory bonus as $\kappa \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}}(\cdot|x) \log \frac{\pi}{\pi_{\text{ref}}}$ where π_{cal} is a fixed calibration distribution. This also falls under GEB by setting $u = -\frac{\pi_{\text{cal}}}{\pi_{\text{ref}}} \log \frac{\pi}{\pi_{\text{ref}}} - \frac{\pi_{\text{cal}}}{\pi_{\text{ref}}} \log \pi_{\text{ref}} + 1$ and f is KL-divergence. The corresponding reparameterized exploratory bonus reduces to $\mathbb{E}_{x \sim \rho, y \sim \pi_{\text{cal}}} - \log \frac{\pi}{\pi_{\text{ref}}} + C(x)$, where $C(x) = \mathbb{E}_{x \sim \rho, y \sim \pi_{\text{ref}}} [-\frac{\pi_{\text{cal}}}{\pi_{\text{ref}}} \log \pi_{\text{ref}} + 1]$. Interestingly, even objectives not explicitly designed for exploration can be reinterpreted through our GEB framework. For instance, Chen et al. [23] augment the DPO loss with an additional term $\kappa \mathbb{E}_{x,y \sim \pi_{\text{ref}}} \sigma(-\beta \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)})$, which was originally introduced to control sample complexity. In our framework, this corresponds to optimistic exploration with $u = -\sigma(-\beta \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)}) + 1$.

5 Experiments

5.1 Experimental Settings

Following prior works [12, 13, 23], we adopt the same iterative online algorithm as in Algorithm 1 with three iterations, aiming to isolate the effects of different exploration bonuses. We adopt two LLM backbones: Llama-3-8B-SFT [17] following prior works, and Mistral-Instruct-v0.3 [26]. The training prompt set is RLHFlow-UltraFeedback [17] as in previous works. URM-LLaMa-3.1-8B [27] serves as the preference oracle. We evaluate the outcome policies on both in-domain and out-of-domain test sets. Specifically, for the in-domain test, we use a held-out test set from UltraFeedback [28], and sample 64 times per prompt with the outcome policy to compare the average reward and win-rates against the base model. We use length-controlled AlpacaEval2 benchmark [29] with GPT-4 as a judge for out-of-domain alignment test, and MATH-500 [30] to evaluate out-of-domain reasoning ability.

Baselines. We adopt f-DPO [25], which extends DPO to the f-divergence class, as the primary baseline. We further compare GEB with three optimistic-exploration methods that incorporate

Table 3: In-domain evaluation on different exploration bonuses. **Boldface** and <u>underline</u> indicate the best and the second-best results, respectively. GEB- π , GEB- $1/\pi$, and GEB- $\arctan(\pi-1)$ corresponds to $1 + \alpha - \pi$, $1/\pi$, and $\arctan(1-\pi) + \alpha$ as in Table 2.

	KL (α=1)		Hel. (α=0.5)		f-KL (α=0)		Avg.		
	WR	AvgR	WR	AvgR	WR	AvgR	WR	AvgR	
Mistral-Instruct-v0.3									
f-DPO	78.42	0.7480	72.69	0.6536	51.11	0.5918	67.40	0.6645	
SELM	77.56	0.7530	-	-	-	-	-	-	
XPO	79.71	0.7492	-	-	-	-	-	-	
VPO	78.57	0.7426	-	-	-	-	-	-	
FEB	78.42	0.7480	71.54	0.6525	47.53	0.5928	65.83	0.6644	
GEB- π	81.00	0.7542	75.48	0.6641	51.68	0.5976	69.39	0.6720	
GEB- $1/\pi$	80.00	0.7554	73.97	0.6541	52.26	0.6051	68.74	0.6715	
GEB-arctanh $(\pi - 1)$	79.71	0.7559	75.69	<u>0.6614</u>	52.76	0.5989	69.39	0.6721	
	LLaMA-3-8B-SFT								
f-DPO	73.11	0.8050	71.11	0.7859	67.38	0.7579	70.53	0.7829	
SELM	74.19	0.8126	-	-	-	-	-	-	
XPO	72.40	0.8119	-	-	-	-	-	-	
VPO	71.61	0.7971	-	-	-	-	-	-	
FEB	73.11	0.8050	68.17	0.7591	67.95	0.7611	69.74	0.7751	
GEB- π	74.34	0.8156	71.68	0.7840	67.67	0.7681	71.23	0.7892	
GEB- $1/\pi$	74.76	0.8102	72.25	0.7859	68.17	0.7591	71.73	0.7851	
GEB-arctanh $(\pi - 1)$	74.98	0.8080	73.26	0.7877	68.89	0.7569	72.38	0.7842	

Table 4: Out-of-domain evaluation on different exploration bonuses with LLaMA-3-8B-SFT. **Bold-face** and <u>underline</u> indicate the best and the second-best results, respectively. GEB- π , GEB- $1/\pi$, and GEB-arctanh($\pi-1$) corresponds to $1+\alpha-\pi$, $1/\pi$, and $\arctanh(1-\pi)+\alpha$ as in Table 2.

	KL(α=1)		Hel.(α=0.5)		f-KL(α=0)		Avg.	
	Alpaca	Math	Alpaca	Math	Alpaca	Math	Alpaca	Math
f-DPO	25.72	67.6	24.73	69.0	17.80	69.2	22.75	68.6
FEB	25.72	67.6	23.75	68.6	19.62	68.6	23.03	68.3
GEB- π	28.27	69.2	25.87	69.6	20.05	71.6	24.73	70.1
GEB- $1/\pi$	<u>26.10</u>	68.4	25.28	70.2	<u>19.80</u>	<u>69.2</u>	23.73	<u>69.3</u>
GEB-arctanh $(\pi - 1)$	24.90	71.0	25.96	67.6	19.62	<u>69.2</u>	23.49	<u>69.3</u>

exploratory bonuses—SELM [12], XPO [13], and VPO [14]. Since the approximations or adaptations in their implementations do not extend beyond the KL divergence, we report their results only under KL. In contrast, we introduce a new baseline, Failed Exploratory Bonus (FEB), which removes these approximations or adaptations, i.e., Eq. 7.

5.2 Results & analyses

GEB delivers robust improvements across different loss designs, divergence classes, and language model backbones. The experimental results are shown in Table 3. Across both backbones, GEB generally outperforms f-DPO and FEB. Under the KL-divergence, GEB displays better or at least on-par performance compared to prior exploratory-bonus methods. Notably, the win-rate increases over 1.82% and 0.94% under the KL-divergence, over 2.36% and 1.29% under the Hellinger Distance, compared with their f-DPO counterpart. GPT-4 evaluation on the Alpaca benchmark also shows consistent performance gains on out-of-domain alignment task. While GEB maintains on par, or usually better results in MATH, showing less performance degradation beyond alignment, known as alignment tax [31, 32].

GEB effectively encourages exploration in small π_{ref} region, yielding more diverse sampling. In Figure 2, we visualize the distribution of $\log \pi_{ref}$ for sampled responses in the last iteration under the KL divergence. When trained with the GEB, the policy model consistently samples more trajectories

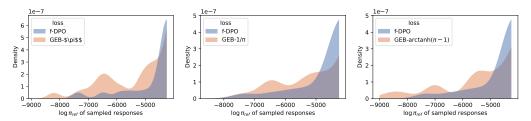


Figure 2: Comparison of $\log \pi_{\text{ref}}$ of sampled response in the last iteration between the general exploratory bonuses and vanilla iterative DPO. GEB- π , GEB- $1/\pi$, and GEB- $\arctan(\pi-1)$ corresponds to $1 + \alpha - \pi$, $1/\pi$, and $\arctan(1-\pi) + \alpha$ as in Table 2

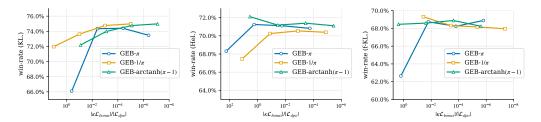


Figure 3: Experiments with different κ . The three graphs are under KL divergence, Hellinger Distance, and forward KL divergence from left to right, respectively. The p, f, tanh in the legends correspond to $1 + \alpha - \pi$, $1/\pi$, $\arctanh(1 - \pi) + \alpha$ in Table 2 respectively.

Table 5: Dist-n of the sampled corpus in the last iteration under the KL divergence.

	dist-1	dist-2	dist-3	dist-4
f-DPO		0.2700		0.8418
$\text{GEB-}\pi$	0.0192	0.2694	0.6323	0.8420
GEB- $1/\pi$	0.0191	0.2738	0.6401	0.8448
GEB-arctanh $(\pi - 1)$	0.0192	0.2730	0.6391	0.8447

with a smaller $\pi_{\rm ref}$ compared to the policy trained by f-DPO loss. This validates our motivation that GEB can encourage sampling trajectories of small $\pi_{\rm ref}$ for optimistic exploration. In Table 5, we further calculate the distinct-n (n=1,2,3,4) for the sampled responses in the last iterations under the KL divergence, which measures the diversity of a corpus. GEB generally has higher diversity scores, validating that GEB incentivizes qualitatively more diverse samples.

The choice of κ . Since the formulation of u in Eq. 10 is flexible, the scale of the GEB term can differ substantially across designs, hence the absolute value of the bonus is less informative. Instead, we examine the relative ratio of the bonus term to the vanilla RL loss $|\kappa \mathcal{L}_{bonus}|/|\mathcal{L}_{RL}|$, which provides a more consistent basis for comparison and offers better practical guidance for tuning κ across diverse settings. As shown in Fig. 3, performance remains stable when the ratio lies within a suitable range (1e-2 to 1e-6 in our case). However, if the ratio is too large, it impedes optimization of the RL objective and degrades performance; if too small, the exploration incentive in uncertain regions diminishes and performance reverts to the vanilla baseline.

6 Conclusion

While recent work proposes exploratory bonuses to operationalize the "optimism in the face of uncertainty" principle, our work shows that the existing theoretical frameworks of exploratory bonuses fail under KL and α -divergence regularization. To address prior theoretical pitfalls, we introduce General Exploratory Bonus (GEB), a novel theoretical framework for sample-efficient RLHF. Our approach provably satisfies the optimism principle and unifies prior heuristic bonuses. We empirically validate GEB on LLM alignment tasks with diverse bonus designs and LLM backbones, showing improved performance and broad applicability across multiple divergence families.

References

- Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao
 Yu, and Yi Wu. Is DPO superior to PPO for LLM alignment? A comprehensive study. In
 Forty-first International Conference on Machine Learning, ICML, 2024.
- Yunhao Tang, Zhaohan Daniel Guo, Zeyu Zheng, Daniele Calandriello, Yuan Cao, Eugene
 Tarassov, Rémi Munos, Bernardo Ávila Pires, Michal Valko, Yong Cheng, and Will Dabney.
 Understanding the performance gap between online and offline alignment algorithms. *CoRR*,
 abs/2405.08448, 2024.
- [3] Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. In *Forty-first International Conference on Machine Learning, ICML*, 2024.
- Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Ramé, Thomas Mesnard, Yao Zhao, Bilal Piot, Johan Ferret, and Mathieu Blondel. Direct language model alignment from online AI feedback. *CoRR*, abs/2402.04792, 2024.
- Isla [5] Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. Iterative reasoning preference optimization. *CoRR*, abs/2404.19733, 2024.
- [6] Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce-style optimization for learning from human feedback in llms. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL*, pages 12248–12267, 2024.
- Nirjhar Das, Souradip Chakraborty, Aldo Pacchiano, and Sayak Ray Chowdhury. Active preference optimization for sample efficient rlhf. In *ICML 2024 Workshop on Theoretical Foundations of Foundation Models*, 2024.
- [8] Kaixuan Ji, Jiafan He, and Quanquan Gu. Reinforcement learning from human feedback with active queries. *CoRR*, abs/2402.09401, 2024.
- [9] Viraj Mehta, Vikramjeet Das, Ojash Neopane, Yijia Dai, Ilija Bogunovic, Jeff G. Schneider, and
 Willie Neiswanger. Sample efficient reinforcement learning from human feedback via active
 exploration. *CoRR*, abs/2312.00267, 2023.
- 133 [10] Han Qi, Haochen Yang, Qiaosheng Zhang, and Zhuoran Yang. Sample-efficient reinforcement learning from human feedback via information-directed sampling. *arXiv preprint* arXiv:2502.05434, 2025.
- 2336 [11] Zichen Liu, Changyu Chen, Chao Du, Wee Sun Lee, and Min Lin. Sample-efficient alignment 337 for Ilms. *arXiv preprint arXiv:2411.01493*, 2024.
- Shenao Zhang, Donghan Yu, Hiteshi Sharma, Ziyi Yang, Shuohang Wang, Hany Hassan, and Zhaoran Wang. Self-exploring language models: Active preference elicitation for online alignment. *CoRR*, abs/2405.19332, 2024.
- Tengyang Xie, Dylan J. Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q*-approximation for sample-efficient RLHF. *CoRR*, abs/2405.21046, 2024.
- Shicong Cen, Jincheng Mei, Katayoon Goshvadi, Hanjun Dai, Tong Yang, Sherry Yang, Dale
 Schuurmans, Yuejie Chi, and Bo Dai. Value-incentivized preference optimization: A unified
 approach to online and offline RLHF. 2025.
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong
 Zhang. Iterative preference learning from human feedback: Bridging theory and practice for
 RLHF under kl-constraint. In Forty-first International Conference on Machine Learning, ICML,
 2024.
- [16] Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong
 Zhang. Iterative preference learning from human feedback: Bridging theory and practice for
 RLHF under kl-constraint. In Forty-first International Conference on Machine Learning, ICML,
 2024.

- Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen
 Sahoo, Caiming Xiong, and Tong Zhang. RLHF workflow: From reward modeling to online
 RLHF. CoRR, abs/2405.07863, 2024.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless
 assistant with reinforcement learning from human feedback. arXiv preprint arXiv:2204.05862,
 2022.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [20] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ah-365 mad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela 366 Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem 367 Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, 368 Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, 369 Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, 370 Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, 371 Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, 372 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab 373 AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco 374 375 Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah 376 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 377 Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason 378 Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya 379 Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, 380 Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Va-381 suden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal 383 Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz 384 Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke 385 de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin 386 Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-387 badur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, 388 Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, 389 Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal 390 Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao 391 Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert 392 Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, 393 Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hos-394 seini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, 395 Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, 396 397 Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, 398 Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal 399 Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, 400 Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin 401 Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, 402 Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine 403 Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, 404 Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, 405 Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay 406 Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit 407 Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, 408 Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, 409 Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, 410 Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, 411 Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, 412

Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, 413 Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester 414 Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon 415 Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, 416 Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin 417 Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, 418 Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, 419 Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank 420 Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, 421 Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan 422 Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison 423 Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, 424 Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, 425 James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff 426 Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, 427 Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh 428 Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun 429 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, 430 Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro 431 Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, 432 Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew 433 Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao 434 Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel 435 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, 436 Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, 437 Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich 438 Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem 439 Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, 440 Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ 443 Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, 444 Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, 445 Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao 446 Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, 447 Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen 448 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, 449 Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, 450 Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim 451 Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, 452 Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu 453 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-454 stable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, 455 Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin 456 Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary 457 DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 458 herd of models, 2024. 459

- 460 [21] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- William Muldrew, Peter Hayes, Mingtian Zhang, and David Barber. Active preference learning
 for large language models. In Forty-first International Conference on Machine Learning, ICML,
 2024.
- 465 [23] Mingyu Chen, Yiding Chen, Wen Sun, and Xuezhou Zhang. Avoiding $\exp(r_{max})$ scaling in rlhf through preference-based exploration. *arXiv preprint arXiv:2502.00666*, 2025.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.

 Advances in Neural Information Processing Systems, 2024.

- [25] Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse KL:
 generalizing direct preference optimization with diverse divergence constraints. In *The Twelfth International Conference on Learning Representations, ICLR*, 2024.
- [26] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,
 Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. CoRR, abs/2310.06825,
 2023.
- 478 [27] Xingzhou Lou, Dong Yan, Wei Shen, Yuzi Yan, Jian Xie, and Junge Zhang. Uncertainty-479 aware reward model: Teaching reward models to know what is unknown. *arXiv preprint* 480 *arXiv:2410.00847*, 2024.
- [28] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan
 Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback.
 In The Thirteenth International Conference on Learning Representations, ICLR, 2024.
- 484 [29] Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *CoRR*, abs/2404.04475, 2024.
- [30] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. arXiv preprint
 arXiv:2305.20050, 2023.
- 489 [31] Michael Noukhovitch, Samuel Lavoie, Florian Strub, and Aaron C. Courville. Language model 490 alignment with elastic reset. In *Advances in Neural Information Processing Systems 36: Annual* 491 *Conference on Neural Information Processing Systems 2023, NeurIPS*, 2023.
- Yong Lin, Hangyu Lin, Wei Xiong, Shizhe Diao, Jianmeng Liu, Jipeng Zhang, Rui Pan,
 Haoxiang Wang, Wenbin Hu, Hanning Zhang, Hanze Dong, Renjie Pi, Han Zhao, Nan Jiang,
 Heng Ji, Yuan Yao, and Tong Zhang. Mitigating the alignment tax of RLHF. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP*, pages
 580–606, 2024.
- [33] Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in
 ML safety. CoRR, abs/2109.13916, 2021.
- [34] Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable
 agent alignment via reward modeling: a research direction. *CoRR*, abs/1811.07871, 2018.
- [35] Xuanming Zhang, Yuxuan Chen, Yiming Zheng, Zhexin Zhang, Yuan Yuan, and Minlie Huang.
 Seeker: Towards exception safety code generation with intermediate language agents framework.
 arXiv preprint arXiv:2412.11713, 2024.
- 504 [36] Xuanming Zhang, Yuxuan Chen, Min-Hsuan Yeh, and Yixuan Li. Metamind: Modeling human social thoughts with metacognitive multi-agent systems. *CoRR*, abs/2505.18943, 2025.
- [37] Min-Hsuan Yeh, Jeffrey Wang, Xuefeng Du, Seongheon Park, Leitian Tao, Shawn Im, and
 Yixuan Li. Position: Challenges and future directions of data-centric ai alignment. In Forty second International Conference on Machine Learning Position Paper Track, 2025.
- [38] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec
 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback.
 Advances in Neural Information Processing Systems, 33:3008–3021, 2020.
- [39] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to
 follow instructions with human feedback. Advances in neural information processing systems,
 35:27730–27744, 2022.
- [40] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland,
 Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning
 from human preferences. In *International Conference on Artificial Intelligence and Statistics*,
 pages 4447–4455. PMLR, 2024.
- [41] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- 522 [42] Vikranth Dwaracherla, Seyed Mohammad Asghari, Botao Hao, and Benjamin Van Roy. Efficient 523 exploration for Ilms. In *Forty-first International Conference on Machine Learning, ICML*, 2024.

- [43] Aobo Kong, Wentao Ma, Shiwan Zhao, Yongbin Li, Yuchuan Wu, Ke Wang, Xiaoqian Liu,
 Qicheng Li, Yong Qin, and Fei Huang. SDPO: segment-level direct preference optimization for
 social agents. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL, pages 12409–12423, 2025.
- Yueqin Yin, Shentao Yang, Yujia Xie, Ziyi Yang, Yuting Sun, Hany Hassan Awadalla, Weizhu Chen, and Mingyuan Zhou. Segmenting text and learning their rewards for improved RLHF in language model. *CoRR*, abs/2501.02790, 2025.
- Yongcheng Zeng, Guoqing Liu, Weiyu Ma, Ning Yang, Haifeng Zhang, and Jun Wang. Token level direct preference optimization. In *Forty-first International Conference on Machine Learn-* ing, ICML, 2024.
- 534 [46] Wendi Li, Wei Wei, Kaihe Xu, Wenfeng Xie, Dangyang Chen, and Yu Cheng. Reinforce-535 ment learning with token-level feedback for controllable text generation. In *Findings of the* 536 *Association for Computational Linguistics: NAACL*, pages 1704–1719, 2024.
- [47] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee,
 Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations, ICLR*, 2024.
- [48] Wendi Li and Yixuan Li. Process reward model with q-value rankings. In *The Thirteenth International Conference on Learning Representations, ICLR*, 2025.
- [49] Lifan Yuan, Wendi Li, Huayu Chen, Ganqu Cui, Ning Ding, Kai Zhang, Bowen Zhou, Zhiyuan
 Liu, and Hao Peng. Free process rewards without process labels. *CoRR*, abs/2412.01981, 2024.
- [50] Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and
 Furu Wei. Reasoning with exploration: An entropy perspective. *CoRR*, abs/2506.14758, 2025.
- [51] Tianyu Zheng, Tianshun Xing, Qingshui Gu, Taoran Liang, Xingwei Qu, Xin Zhou, Yizhi Li,
 Zhoufutu Wen, Chenghua Lin, Wenhao Huang, Qian Liu, Ge Zhang, and Zejun Ma. First return,
 entropy-eliciting explore. *CoRR*, abs/2507.07017, 2025.
- Fang Wu, Weihao Xuan, Ximing Lu, Zaïd Harchaoui, and Yejin Choi. The invisible leash: Why RLVR may not escape its origin. *CoRR*, abs/2507.14843, 2025.
- Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui
 Chen, Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji
 Song, Bowen Yu, Gao Huang, and Junyang Lin. Beyond the 80/20 rule: High-entropy minority
 tokens drive effective reinforcement learning for LLM reasoning. CoRR, abs/2506.01939, 2025.
- [54] Han Zhong, Wei Xiong, Sirui Zheng, Liwei Wang, Zhaoran Wang, Zhuoran Yang, and Tong
 Zhang. GEC: A unified framework for interactive decision making in mdp, pomdp, and beyond.
 CORR, abs/2211.01962, 2022.

A Related Works

Alignment & RLHF. Alignment [33, 34, 35, 36, 37] aims to ensure AI systems act in accordance with human values, preferences, and goals; and it has become a critical field in AI research. To steer language models to match human preferences, reinforcement Learning from Human Feedback (RLHF) [38, 39] acheives great success and has become the standard alignment pipeline. However, its computational complexity has motivated a family of Direct Preference Optimization (DPO) [24, 40, 41] that forgo explicit reward modeling. Despite their efficiency, recent researchers [3, 1, 16] reemphasize the significance of online sampling.

Optimistic exploration of RLHF. To address the computational overheads of passive exploration in RLHF, which samples trajectories just based on randomness, some existing attempts have been devoted to sample-efficient RL algorithms. Most of works [7, 8, 42, 9, 22] adhere to the principle of optimism, proposing specialized prompt or response selection strategies to emphasize uncertain samples. While some research [11, 27] propose uncertainty-aware reward models with epistemic neural networks or bootstrap ensembles, these methods introduce additional cost. Some research also addresses the sample efficiency with different theoretical foundations, such as information theory [10], preference-incentive exploration [23]. Notably, several works [12, 13, 14] introduce different exploratory bonuses, which can implement optimism toward uncertainty without additional computes. However, they only focus on KL-divergence and their theoretical framework cannot result in real optimism as shown in Section 3.2.

Efficient RL for LLM. Beyond optimistic exploration, some research proposes fine-grained signals for RL learning. For instance, several research propose segment-level [43, 44] or token-level [45, 46] reward function for alignment or text control. Notably, for reasoning tasks, process reward model [47, 48, 49] which provides step-wise feedback for solutions has shown promise effectiveness. On the other hand, recent research [50, 51, 52] on LLM reasoning reveal that high-entropy tokens guide the model toward diverse reasoning paths. Training with only high-entropy tokens are more beneficial for reasoning performance [53]. While our approach is highly extensible, we believe the orthogonal methods can be further incorporated with our general exploratory bonus.

B Optimism Failure of previous works

B.1 Extension beyond α -divergence

The following theorem formally proves that the exploratory bonus $-\kappa \max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$ cannot encourage optimism for more general divergence class.

Theorem 3.3 When f belongs to f-divergence, and the reward function is obtained by $\hat{r} = \arg\min_r [\mathcal{L}_{BT}(\mathcal{D}_t, r) - \kappa \max_{\pi} \mathcal{J}_{\beta,f}(\pi, r)]$ and the policy is updated by $\arg\max_{\pi} \mathcal{J}_{\beta,f}(\pi, r_t)$, the bonus term $-\kappa \max_{\pi} \mathcal{J}_{\beta,f}(\pi, r)$ induces the policy model π to coincide with π_{ref} when xf''(x) is a monotone function.

Proof By Lemma. 3.2, we reparameterize the bonus term for optimistic reward-modeling to Eq. 8. Denote h(u) = uf'(u) - f(u). For a fixed prompt x, we formulate the training process as a constrained problem as follows,

$$\arg\max_{\pi} \mathbb{E}_{y \sim \pi_{\text{ref}}(\cdot|x)} h(\frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)}) \quad s.t. \quad \sum_{y} \pi(y|x) = 1 \quad \text{and} \quad \forall y, \pi(y|x) > 0. \tag{13}$$

Then we can apply the Lagrange multiplier as

$$\mathcal{L} = \mathbb{E}_{y \sim \pi_{\text{ref}}} h\left(\frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)}\right) - \mu\left(\sum_{y} \pi(y|x) - 1\right) - \sum_{y} \eta(y)\pi(y|x),\tag{14}$$

where μ , η are the dual variables. Then we utilize the Karush-Kuhn-Tucker (KKT) conditions for the given optimization problem. The complementary slackness requires that $\forall y, \eta(y)\pi(y|x) = 0$. The stationary condition requires

$$\frac{\partial \mathcal{L}}{\partial \pi(y|x)} = h'(\frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)}) - \mu - \eta(y) = 0.$$
 (15)

We denote $S_y = \{y | \pi(y|x) > 0\}$, we have $\forall y \in S_y, \eta(y) = 0$. Since h'(u) = uf''(u) > 0 due to the convexity of $f(\cdot)$, we can obtain $\forall y \in S_y, \frac{\pi(y|x)}{\pi_{\rm ref}(y|x)}$ is a constant. Then applying the normalisation 600 601 constraint, $\mathbb{E}_{y \sim \pi_{\text{ref}}(\cdot|x)} \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} = 1$. Hence, the unique interior optimum is $\pi^*(y|x) = \pi_{\text{ref}}(y|x)$. \square 602

The theorem implies that the reparameterized exploratory bonus attains its maximum only when π 603 and π_{ref} coincide. The condition that xf''(x) is a monotone function is satisfied by α -divergence 604 and beyond, e.g. Pearson χ^2 . Hence, the exploratory bonus $-\kappa \max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$ in the reward 605 training objective generally contradicts the optimism, since it cannot encourage trajectories with 606 607 small initialized possibility.

B.2 Prior adaptions & approximations cannot generalize 608

Though the theoretical framework of prior exploratory bonus fails, their empirically implemented 609 loss remains effective through different adaptions and approximations. However, in this subsection, 610 we show these adaptions and approximations are inextensible beyond KL-divergence class. 611

Zhang et al. [12] adapt the formulation of $\mathcal{J}_{\beta,f}(\pi,r)$ in $-\kappa \max_{\pi} \mathcal{J}_{\beta,f}(\pi,r)$ as 612

$$\mathcal{J}'_{\beta,f}(\pi,r) = E_{x,y \sim \pi, y' \sim \pi_{\text{ref}}}[r(x,y) - r(x,y')] - \beta D_{KL}(\pi|\pi_{\text{ref}})$$

$$\tag{16}$$

be zero after re-parameterization as shown in Lemma 3.1, thus the sole reparameterized bias term will remain as $-\mathcal{E}_{y'\sim\pi_{\text{ref}}}\log\pi(y'|x)$. Since $\mathcal{J}_{\beta,f}(\pi,r)$ cannot be reparameterized to zero except 615 KL-divergence, this adaption cannot generalize then. 616 In the derivations of Cen et al. [14], it utilizes an ideal distribution π_{cal} which should satisfy 617 $\mathbb{E}_{y \sim \pi_{cal}} r(x, y) = 0$. Since π_{cal} is practically unobtainable, it uses the rejected responses to approximate $\mathbb{E}_{\pi_{cal}}$, which does not satisfy the predefined condition of π_{cal} thus not rigorously coherent 619 to the theory. While the regret decomposition of Xie et al. [13] relies on the logarithm form of 620 KL-divergence, thus inextensible to broader divergence class. 621 In contrast, our general exploratory bonus can seamlessly incorporate iterative online RLHF algorithm 622 and can naturally extend to the entire α -divergence class. All prior bonuses mentioned above can be 623

which adds a bias in reward expectation term. Under KL-divergence, the original $\mathcal{J}_{\beta,f}(\pi,r)$ will

Regret Bound 625

613

614

In derivations, we utilize the theoretical tools in [12, 13, 14]. First, we make some standard statistical 626 assumptions following Cen et al. [14]. 627

Assumption C.1 For a reward function r, and a random function $R(\cdot)$, and any trajectory τ , we 628 have $-R_{max} < r(\tau), R(r(\tau)) < R_{max}$, where R_{max} is a constant. 629

This is an assumption generally made for theoretical analyses of RLHF. Note that R_{max} is measurable 630 and controllable in practice. Then we introduce the assumption of the reward class proposed in Cen 631 et al. [14], which offers a regularization mechanism to incorporate additional policy preferences in 632 the subsequent derivations. 633

Assumption C.2 We assume that $r^* \in \mathcal{R}$, where 634

encompassed by our theoretical framework.

$$\mathcal{R} = \{ r : \mathbb{E}_{x \sim \rho, \tau \sim \pi_{cal}(\cdot|x)} r(x, y) = 0 \}, \tag{17}$$

where ρ is the prompt distribution and π_{cal} is a fixed calibration distribution independent of the 635 algorithm. 636

We also introduce the preference generalized eluder coefficient proposed in Zhang et al. [12], an 637 extension of the generalized eluder coefficient [54], which connects prediction error and in-sample estimation error.

Definition C.1 Let $f_r(x,y,y') = r_t(x,y) - r^*(x,y)$. For a reward function class \mathcal{R} , we define the preference generalized eluder coefficient as the smallest d_{PGEC} as

$$\sum_{t=1}^{T} \mathbb{E}_{x \sim \rho, y \sim \pi^{t}(\cdot | x), y' \sim \pi_{cal}} [f_{r_{t}}(x, y, y') - f_{r^{*}}(x, y, y')]$$

$$\leq \sqrt{d_{PGEC} \sum_{t=1}^{T} \mathbb{E}_{x \sim \rho, y \sim \widetilde{\pi}_{t}(\cdot | x), y' \sim \pi_{cal}} [f_{r_{t}}(x, y, y') - f_{r^{*}}(x, y, y')]^{2} + 4\sqrt{d_{PGEC}T}}, \quad (18)$$

With the above assumptions and the theoretical tool, we can have the following regret boundary. 643

- **Theorem C.1** Let $\mathcal{J}_{\beta,f}(\pi,r) = \mathbb{E}_{x \sim \rho, y \sim \pi(\cdot|x)} r(x,y) \beta \mathbb{E}_{x \sim \rho, y \sim \pi_{ref}(\cdot|x)} f(\frac{\pi(y|x)}{\pi_{ref}(y|x)})$, when the hyperparameter of the loss Eq. 3 $\kappa = \sqrt{\frac{\log(T|\mathcal{R}|\delta^{-1})}{(\gamma d_{\text{PGEC}}T)}} (32R_{max}e^{4R_{max}})^{-1}$, with probability at least 1δ , the regret can be bounded as follows,

642

$$\sum_{t=1}^{T} \mathcal{J}_{\beta,f}(\pi^*, r^*) - \mathcal{J}_{\beta,f}(\pi^t, r^*) \le \mathcal{O}(R_{max} e^{4R_{max}} T \sqrt{d_{PGEC} \gamma \log(|\mathcal{R}|\delta^{-1})}), \tag{19}$$

- where $\gamma = \sup_{x,y} \frac{\pi}{\pi_{cal}}$, r^* and π^* are ground-truth reward function and corresponding optimal policy with $\pi^* = \arg \max_{\pi} \mathcal{J}_{\beta,f}(\pi, r^*)$.
- *Proof* First, we can decompose the regret function as in Cen et al. [14] as follows,

$$\underbrace{\sum_{t=1}^{T} [\mathcal{J}_{\beta,f}(\pi^*, r^*) - \mathcal{J}_{\beta,f}(\pi^t, r^t)]}_{\text{Term 1}} + \underbrace{\sum_{t=1}^{T} [\mathcal{J}_{\beta,f}(\pi^t, r^t) - \mathcal{J}_{\beta,f}(\pi^t, r^*)]}_{\text{Term 2}}.$$
 (20)

- Then, we will bound term 1 and term 2 individually and combine them at last.
- **Bound term 1.** First, we connect the term 1 with $\max_{\pi} \mathcal{J}_{\beta,f}(\pi, R(r_t)) \max_{\pi} \mathcal{J}_{\beta,f}(\pi, R(r^*))$. 651
- When π^* is the optimal π for $\max_{\pi} J_{\beta,f}(\pi, r^*)$, we have 652

Term
$$1 \le J_{\beta,f}(\pi^*, r^*) - J_{\beta,f}(\pi^*, r_t) \le \sup_{x,y} \frac{\pi^*}{\pi_t} \mathbb{E}_{x \sim \rho, y \sim \pi_t}(r^* - r^t).$$
 (21)

Similarly, we can obtain its lower bound as $-2R_{max}$. Then, we have

Term
$$1 \le \sum_{t=1}^{T} \left[\max_{\pi} \mathcal{J}_{\beta,f}(\pi, R(r_t)) - \max_{\pi} \mathcal{J}_{\beta,f}(\pi, R(r^*)) \right] + 4R_{max}T.$$
 (22)

Bound term 2. First, we utilize the preference generalized eluder coefficient to connect the prediction error to in-sample error.

Term 2 =
$$\sum_{t=1}^{T} \mathbb{E}_{x \sim \rho, y \sim \pi^{t}(\cdot|x), y' \sim \pi_{cal}} [f_{r_{t}}(x, y, y') - f_{r^{*}}(x, y, y')]$$
 (23)

$$\leq \frac{\eta T d_{\text{PGEC}}}{4} + 4\sqrt{d_{\text{PGEC}}T} + \frac{1}{\eta} \sum_{t=1}^{T} \mathbb{E}_{x \sim \rho, y \sim \tilde{\pi}_{t}(\cdot|x), y' \sim \pi_{cal}} [f_{r_{t}}(x, y, y') - f_{r^{*}}(x, y, y')]^{2}, \quad (24)$$

- where the first equality uses the property of the reward class in Assumption C.2, and the inequality follows Definition C.1 with Cauchy-Schwarz inequality. Then, we bound the squared in-sample error
 - $\mathbb{E}_{x \sim \rho, \tau \sim \widetilde{\pi}_{t}} [f_{r_{t}}(x, y, y') f_{r^{*}}(x, y, y')]^{2} \leq \gamma \mathbb{E}_{x \sim \rho, \tau, \tau' \sim \widetilde{\pi}_{t}} [f_{r_{t}}(x, y, y') f_{r^{*}}(x, y, y')]^{2}$ (25)

$$\leq \gamma (32R_{max}e^{4R_{max}})^2 \mathbb{E}_{x \sim \rho, \tau, \tau' \sim \widetilde{\pi}_t} [\sigma(f_{r_t}(x, y, y')) - \sigma(f_{r^*}(x, y, y'))]^2 \tag{26}$$

$$\leq 8\gamma (32R_{max}e^{4R_{max}})^2 \mathbb{E}_{x\sim\rho,\tau,\tau'\sim\widetilde{\pi}_t(\cdot|x)} D_H^2(P_{r_t}(\cdot|\tau,\tau') \|P_{r^*}(\cdot|\tau,\tau')), \tag{27}$$

where $\gamma=\sup_{x,y}\frac{\widetilde{\pi}_t}{\pi_{cal}}$, and the second inequality utilizes the Lemma C.8 in Xie et al. [13], the third inequality uses $(x-y)^2<4(x+y)(\sqrt{x}-\sqrt{y})$. Refer to the Lemma C.6 in Xie et al. [13], we have

$$\sum_{i < t} \mathbb{E}_{x \sim \rho, \tau, \tau' \sim \widetilde{\pi}_t(\cdot|x)} D_H^2(P_{r_t}(\cdot|\tau, \tau') \| P_{r^*}(\cdot|\tau, \tau')) \le L_{BT}^{(t)}(r_t) - L_{BT}^{(t)}(r^*) + 2\log(|\mathcal{R}|\delta^{-1})$$
 (28)

where $L_{BT}^{(t)}(r) = \sum_{i < t} \mathbb{E}_{y,y' \sim \mathcal{D}_{\square}} - \log \sigma(f_r(x,y,y'))$ is the vanilla BT loss for reward modeling. Finally, the term 2 can be bounded by

Term
$$2 \le 4\sqrt{d_{\text{PGEC}}T} + \frac{\eta T d_{\text{PGEC}}}{4} + \frac{8\gamma}{\eta} (32R_{max}e^{4R_{max}})^2 (L_{BT}^{(t)}(r_t) - L_{BT}^{(t)}(r^*) + 2T\log(|\mathcal{R}|\delta^{-1})).$$
 (29)

Bound the regret. Since the r_t is optimized by $L_{BT}^{(t)}(r_t) - \sum_{i=1}^T \kappa \max_{\pi} \mathcal{J}_{\beta,f}(R(r_t),\pi)$ we have $r_t = \arg\min_{r \in \mathcal{R}} L_{BT}^{(t)}(r_t) - \sum_{i=1}^T \kappa \max_{\pi} \mathcal{J}_{\beta,f}(R(r_t), \pi)$ Therefore, when $\eta = \frac{1}{2} \sum_{i=1}^T \kappa \max_{\pi} \mathcal{J}_{\beta,f}(R(r_t), \pi)$ $4\sqrt{\frac{2\gamma\log(T|\mathcal{R}|\delta^{-1})}{Td_{\mathrm{PGEC}}}}(32R_{max}e^{4R_{max}})$ and $\kappa=\frac{\eta}{4\gamma}(32R_{max}e^{4R_{max}})^{-2}$, the regret can be bounded

$$\operatorname{Regret} \le 4\sqrt{d_{\mathrm{PGEC}}T} + \sqrt{2^{3}\gamma d_{\mathrm{PGEC}}\log(|\mathcal{R}|\delta^{-1})}(32R_{max}e^{4R_{max}}T) + 4R_{max}T. \tag{30}$$

667

Experiments D 668

669

670

671

672

673

674

675

676

678

679

680

681

682

683

684

Implementation Details

Algorithm. Following prior works Zhang et al. [12], Xie et al. [13], Chen et al. [23], we adopt the same algorithmic backbone for empirical validation to explicitly show the effect of different exploratory bonus in loss function. This algorithm bypasses the reward modeling in each iteration through reward reparameterization, known as iterative DPO [17]. Previous works further reparameterizes the bonus term to incorporate the algorithm. Since iterative DPO can seamlessly extend to f-divergence, we also follow prior works to reparameterize our general exploratory bonus. The detailed algorithm can be formulated as in Algorithm 1.

Algorithm 1 Iterative Online Algorithm with Exploratory Bonus

Input: Reference model π_{ref} , iteration number T, prompt set for each interation $\mathcal{D}_1, \dots, \mathcal{D}_T$, reward function r;

```
Output: Trained model \pi_T;
 1: for iteration t = 1, 2, ..., T do
```

for $x \in \mathcal{D}_t$ do

3: $y_1, y_2 \sim \pi_{\text{ref}}(\cdot|x)$ and obtain the rewards $r(y_1), r(y_2)$;

Rank the reward and denote y^+, y^- as the preferred and dispreferred response between y_1, y_2 and update $D_t = \{x, y^w, y^l\}$;

5: $\pi_t = \arg\min_{\pi} \mathcal{L}_{DPO} - \kappa \mathcal{L}_{bonus}(\pi)$

6: update π_{ref} with π_t (optional)

7: end for

8: end for

Hyperparameter settings and environments. All experiments are conducted on two NVIDIA H200 GPUs. When training and sampling, the max length is set to 2048. For training, the batch size per device is set to 2; we enable the gradient checkpointing and the gradient accumulation step is set to 64; the learning rate is 5e-7 with cosine scheduler, and the warm up ratio is 0.03. In main experiments, we use the best performance with κ with a suitable ratio range to f-dpo loss across 1, 1e-2, 1e-4, 1e-6, 1e-8. For sampling, the temperature is set to 1. For in-domain evaluation and MATH evaluation, we set temperature to 0.6 and top-p to 0.9; we use the default setting of alpaca-eval.

E Proofs of Lemma 4.1

Lemma 4.1 Denote two objectives as $h(u) = \mathbb{E}_{x \sim \rho, y \sim \pi_{ref}} u f'(u) - f(u)$ and $\hat{h}(u) = \mathbb{E}_{x \sim \rho, y \sim \pi_{ref}} \frac{u}{Z(x)} f'(u) - f(\frac{u}{Z(x)})$, where u varies with π and π_{ref} . When $[f'(u) + u f''(u) - f(\frac{u}{Z(x)})]/[Z(x)uf''(u)] = \Lambda(x)$ is constant in y and $\Lambda(x) > 0$, minimizing the two objectives $\min_{\pi} -h(u)$ and $\min_{\pi} -\hat{h}(u)$ induce the same class of policies.

Proof For a fixed x, with the Lagrange multiplier and KKT conditions, when $\pi(y|x)>0$, we 691 obtain

$$\frac{\partial h}{\partial \pi(y|x)} = u(x,y)f''(u(x,y)) \cdot \frac{\partial \pi_{\text{ref}}(y|x)u(x,y)}{\partial \pi(y|x)} = \mu_1(x), \tag{31}$$

where $\mu_1(x)$ is a dual variable with respective to x. When $\frac{f'(u)+uf''(u)-f(\frac{u}{Z})}{Zuf''(u)}=\Lambda(x)$, we have

$$\mu_1(x)\Lambda(x) = \frac{1}{Z}(f'(u) + uf''(u) - f'(\frac{u}{Z})) \cdot \frac{\partial \pi_{\text{ref}}(y|x)u(x,y)}{\partial \pi(y|x)},\tag{32}$$

which directly follows by the KKT conditions of $\hat{h}(x)$, i.e. $\frac{\partial \hat{h}}{\partial \pi(y|x)} = \mu_2(x)$ where $\mu_2(x)$ is a dual variable equals to $\mu_1(x)\Lambda(x)$. Hence, every policy that satisfies the stationary condition for h also satisfies it for \hat{h} . Since $\Lambda(x)>0$, the second-order derivative $\frac{\partial h}{\partial^2 \pi(y|x)}$ and $\frac{\partial \hat{h}}{\partial^2 \pi(y|x)}$ have the same sign, which indicates they share the same local minima. Hence, minimizing the two objectives $\min_{\pi} -h(u)$ and $\min_{\pi} -\hat{h}(u)$ induce the same class of policies.

699 F Statement & Limitations

We have included all implementation details, hyperparameters, and training procedures in the paper and appendix. Our code and scripts for reproducing the experiments are available through the anonymous GitHub repository to obey the double-blind policy, and will be further made publicly available upon publication.

This work studies reinforcement learning from human feedback (RLHF) using only publicly available or synthetic data, without new human subject collection. Here, by providing a rigorous theoretical framework with strong empirical evidence, we pursue a high standard of scientific excellence. We also take into account inclusiveness to make all our visualizations accessible to the unprivileged group of people, by producing figures distinguished by light, shade, and marker. While RLHF has the potential to amplify biases or harmful behaviors if misused, our work is intended solely to advance safe and responsible research, and we encourage its application in alignment with ethical standards.

GEB is generally based on online iterative RLHF. This online RLHF backbone has an off-policy instinct. We do not explore whether our GEB can also be seamlessly incorporated into more on-policy algorithms. Meanwhile, our experiments only focuses on the alignment task; whether GEB can benefit more general task is still mysterious. Nonetheless, our paper shows the failure of the existing theoretical frameworks of exploratory bonuses, and introduce General Exploratory Bonus (GEB), a novel theoretical framework for sample-efficient RLHF. Our approach provably satisfies the optimism principle and unifies prior heuristic bonuses. The empirical results also show improved performance and broad applicability across multiple divergence families.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The failure of existing theoretical framework is demonstrated in §3.2. The theoretical framework is introduced in §4. The empirical studies are in §5

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The limitations are discussed in §F.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
 only tested on a few datasets or with a few runs. In general, empirical results often
 depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We generally provide proofs directly after the lemma or theorem, while some of proofs are supplemented in Appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The implementation details can be found in §5 and §D.1. Moreover, we provide the reproducible code and scripts here.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
 well by the reviewers: Making the paper reproducible is important, regardless of
 whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The datasets we use are all open-sourced, and we have provided the reproducible code and scripts here..

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
 including code, unless this is central to the contribution (e.g., for a new open-source
 benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
 proposed method and baselines. If only a subset of experiments are reproducible, they
 should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The implementation details can be found in §5 and §D.1. Moreover, we provide the reproducible code and scripts here.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: The repeated training on large-scale LLMs are computation-costly, but we use multiple variants of GEB and experiments of different hyper-parameters to validate the stability of GEB.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891 892

894

895

896

897

898

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

925

Justification: The implemented details are in §D.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research is conducted according to the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]
Justification: In §F.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal
 impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [No]

Justification: A aim of RLHF is to reduce the occurrence of jailbreaking behaviors. Existing safeguard strategies can be generally applied to our outcome policies.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The creators or original owners of assets (e.g., code, data, models), used in the paper, are properly credited.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018 1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The usage of outcome policies follows the standard usage of LLMs in the huggingface package.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.