000 **REWARD-AUGMENTED DATA ENHANCES** DIRECT 001 Preference Alignment of LLMs 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Preference alignment in Large Language Models (LLMs) has significantly improved their ability to adhere to human instructions and intentions. However, existing direct alignment algorithms primarily focus on relative preferences and often overlook the qualitative aspects of responses, despite having access to preference data that includes reward scores from judge models during AI feedback. Striving to maximize the implicit reward gap between the chosen and the slightly 016 inferior rejected responses can cause overfitting and unnecessary unlearning of the high-quality rejected responses. The unawareness of the reward scores also drives the LLM to indiscriminately favor the low-quality chosen responses and fail to generalize to responses with the highest rewards, which are sparse in data. To overcome these shortcomings, our study introduces reward-conditioned LLM policies that discern and learn from the entire spectrum of response quality within the dataset, helping extrapolate to more optimal regions. We propose an effective vet simple data relabeling method that conditions the preference pairs on quality scores to construct a reward-augmented dataset. This dataset is easily integrated with existing direct alignment algorithms and is applicable to any preference dataset. The experimental results across instruction-following benchmarks including AlpacaEval 2.0, MT-Bench, and Arena-Hard-Auto demonstrate that our approach consistently boosts the performance of DPO by a considerable margin across diverse models such as Zephyr, Mistral, Qwen2, Llama3.1, Gemma2, and SPPO. Additionally, on six academic benchmarks including GSM8K, GPQA, MUSR, TruthfulQA, BBH, and ARC, our method improves their average accuracy. When applying our method to on-policy data, the resulting DPO model outperforms various baselines and achieves state-of-the-art results on AlpacaEval 2.0. Through comprehensive ablation studies, we demonstrate that our method not only maximizes the utility of preference data but also mitigates the issue of unlearning, demonstrating its broad effectiveness beyond mere dataset expansion.

004

010 011

012

013

014

015

017

018

019

021

023

025

026

027

028

029

031

032

034

1 INTRODUCTION

039 Reinforcement Learning from Human Feedback (RLHF) has recently seen remarkable success in 040 aligning Large Language Models (LLMs) to follow instructions with human intentions. In this 041 approach, AI-generated feedback serves as a stand-in for human preferences, assessing and ranking 042 responses to prompts to construct a preference dataset. This dataset is then utilized in preference 043 optimization algorithms to fine-tune LLMs. Among them, direct preference alignment (Rafailov 044 et al., 2024b; Azar et al., 2023; Zhao et al., 2023; Ethayarajh et al., 2024) that bypasses the need for an explicit reward model has garnered interest for their simplicity and cost efficiency. However, these algorithms mainly concern relative preferences and often overlook the quality of responses 046 and their gaps, leading to limitations in their effectiveness. 047

048 Specifically, direct alignment algorithms such as DPO (Rafailov et al., 2024b) focus on maximizing the implicit reward difference between accepted and rejected responses. This approach can lead to overfitting, as high-quality but rejected responses are unnecessarily unlearned (Adler et al., 2024). 051 Even worse, since the dataset provides only a sample estimate of true preferences, the rejected responses can actually be more aligned with human preferences than the accepted ones in expectation. 052 Similarly, due to the unawareness of the responses' qualities, direct alignment will also result in the indiscriminate learning of the chosen responses, even when they are of low quality. As a result, the directly aligned LLMs often struggle to differentiate between responses of varying quality and
 fail to generalize effectively to more optimal or the highest-reward responses that are sparse in the
 preference data, which is another limitation.

057 To address these issues, we propose learning reward-conditioned policies as a straightforward fix 058 to the above issues. By optimizing the LLM to generate responses conditioning on their quali-059 ties, the model is allowed to discern and leverage patterns within responses of varied quality. As 060 a result, learning from both chosen and rejected responses alleviates the issue of unnecessarily un-061 learning high-quality rejected responses; distinguishing between varying-quality chosen responses 062 alleviates the issue of indiscriminately accepting low-quality ones. By identifying common patterns 063 in responses of similar quality and distinguishing them from those of differing quality, the LLM 064 becomes more adept at generalizing to more optimal responses that are sparse in data.

065 With this motivation, we introduce an effective yet simple data relabeling method to construct 066 reward-augmented datasets. We define a goal-conditioned reward using an indicator function that 067 compares the goal reward with the actual quality score, such as the reward value given by the judge 068 model during AI feedback. This allows us to relabel each preference pair, generating two new pairs 069 conditioned on the reward goals of both the chosen and rejected responses. The resulting augmented 070 dataset, which contains these newly conditioned pairs, can enhance the performance of existing di-071 rect alignment algorithms. Our method can be applied to any preference dataset and followed by off-the-shelf direct alignment algorithms to boost their performance. 072

073 In experiments, we first apply our method on UltraFeedback (Cui et al., 2023) and perform DPO 074 (Rafailov et al., 2024b) on this reward-augmented preference dataset by fine-tuning on various mod-075 els, including Zephyr-7B- β (Tunstall et al., 2023b), Mistral-7B-Instruct-v0.3 (Jiang et al., 2023a), 076 Qwen2-7B-Instruct (Yang et al., 2024), Llama-3.1-8B-Instruct (Dubey et al., 2024), Gemma-2-9B-077 It (Team et al., 2024), and SPPO (Wu et al., 2024). The results show that our method consistently boosts the performance of these models as well as their DPO models by a large margin on instruction-following benchmarks such as AlpacaEval 2.0 (Dubois et al., 2024), MT-Bench (Zheng 079 et al., 2024), and Arena-Hard-Auto (Li et al., 2024b). Our method also improves the average accuracy on a variety of academic benchmarks (GSM8K, GPQA, MUSR, TruthfulQA, BBH, and ARC). 081 Moreover, our findings also demonstrate an improved utility of the preference data: a subsequent 082 round of DPO using the reward-augmented data can still significantly enhance the model fine-tuned 083 with DPO; relabeling the binarized preference dataset with the DPO implicit reward leads to further 084 performance gains. Additional ablation studies also suggest that our method addresses the problem 085 of unlearning and is superior not just due to the increased dataset size. When applied to on-policy data, our method enhances the DPO model, enabling it to surpass various baselines and achieve 087 state-of-the-art performance on AlpacaEval 2.0.

880

090 091

092

093

094

095 096

2 BACKGROUND

Consider a language model π ∈ Δ^X_Y that takes the prompt x ∈ X as input and outputs the response y ∈ Y, where X and Y are spaces of prompts and responses, respectively. Given the prompt x ∈ X, a discrete probability distribution π(· | x) ∈ Δ_Y is generated, where Δ_Y is the set of discrete distributions over Y. We define the true human preference distribution as

$$p^*(y_1 \succ y_2 \mid x) := \mathbb{E}_h [\mathbb{1}(h \text{ prefers } y_1 \text{ over } y_2 \text{ given } x)],$$

where h denotes the human rater and the expectation is over h to account for the randomness of the human raters' choices. After pretraining and Supervised Fine-Tuning (SFT), Reinforcement Learning from Human or AI Feedback (Ouyang et al., 2022; Bai et al., 2022b) is typically employed to enhance the ability of the language model to follow instructions with human preferences.

101

RL from AI Feedback (RLAIF). The RLAIF framework involves two major steps: preference dataset construction with AI feedback and preference optimization. As a surrogate for human preference, AI feedback, including LLM-as-Judge (Zheng et al., 2024; Cui et al., 2023) and Reward-Model-as-Judge (Adler et al., 2024; Dong et al., 2024), can be used to rank responses and generate preference pairs. Specifically, consider the judge model $r(x, y) : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ that outputs a scalar reward value representing the quality of y under x. For each prompt $x \in \mathcal{X}$, two responses, y_1 and y_2 , are independently sampled—either from the same reference model (Xiong et al., 2024; Wu

123

124 125

126

127 128

129 130

131

146

157

158

161

et al., 2024) or several different models (Zhu et al., 2023a; Zhang et al., 2024). Then $r(x, y_1)$ and $r(x, y_2)$ are evaluated to determine the preferred response $y_w = \operatorname{argmax}_{y \in \{y_1, y_2\}} r(x, y)$ and dispreferred response $y_l = \operatorname{argmin}_{y \in \{y_1, y_2\}} r(x, y)$. By sampling responses and ranking them for a set of N prompts, we get a preference dataset: $\mathcal{D}_N = \{(x^i, y^i_w, y^i_l)\}_{i=1}^N$. For the simplicity of our discussions, we assume that the reward function r is bounded in $[0, r_{\text{max}}]$.

114 **Direct Alignment from Preference.** The objective for the LLM $\pi \in \Delta_{\mathcal{Y}}^{\mathcal{X}}$ is to maximize the KLregularized expected reward. Recent works (Azar et al., 2023; Zhao et al., 2023; Tunstall et al., 2023b; Ethayarajh et al., 2024) proposed to align the LLM directly with the preference data by deriving the preference loss as a function of the LLM by the change of variables. Among them, the Direct Preference Optimization (DPO) (Rafailov et al., 2024b) loss has the following form:

$$\mathcal{L}_{\text{DPO}}(\pi; \mathcal{D}_N) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}_N} \left[\log \sigma \left(\beta \log \frac{\pi(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right]$$

where β is a hyperparameter corresponding to the KL divergence regularization, $\sigma(\cdot)$ is the logistic function, and π_{ref} is some reference LLM policy, such as the SFT model.

3 REWARD-CONDITIONING ADDRESSES LIMITATIONS OF DIRECT PREFERENCE ALIGNMENT

3.1 LIMITATIONS OF DIRECT ALIGNMENT FROM PREFERENCE

We will first demonstrate the limitations of vanilla direct alignment over the preference data.

132 High-Quality Rejected Responses are Unnecessarily Suppressed. The dataset \mathcal{D}_N often con-133 tains preference pairs where the rejected response y_l is only marginally worse than the chosen one 134 y_w . Direct alignment algorithms, however, primarily focus on relative preferences and are unaware 135 of the responses' quality values and gaps. Striving to maximize the reparameterized reward gap 136 between the chosen and rejected responses will risk overfitting and unnecessary "unlearning", i.e., 137 probability decrease, of high-quality responses, potentially diminishing the model's performance 138 by discarding valuable alternatives. Furthermore, in such a finite data regime where only a sample 139 estimate of the true preference is accessible, it can be very possible that $p^*(y_l \succ y_w \mid x) > 0.5$, i.e., y_l is in fact more preferred than y_w in expectation. This issue becomes even more pronounced when 140 the preference data generated with the imperfect judge model is noisy. 141

We illustrate this limitation with the example in Table 1, where we define the maximum reward r_{max} as 10. For $\mathcal{D}_{N=1}$ that contains a single preference pair¹ with reward $r(x, y_1) = 9$ and $r(x, y_2) = 8$, the optimal policy learned from $\mathcal{D}_{N=1}$ is $\pi^*(y_1 \mid x) = 1$. This causes the model to avoid generating y_2 , a response of nearly equivalent quality.

Low-Quality Chosen Responses are Indiscriminately Learned. For a similar reason, direct alignment algorithms also indiscriminately reinforce the chosen responses. As illustrated in Table 2, when $\mathcal{D}_{N=2}$ contains two preference pairs, where one of the chosen responses, y_2 , is of low quality, π^* still indiscriminately generates y_2 with an arbitrary probability $0 \le a \le 1$, i.e., $\pi^*(y_2 \mid x) = a$.

Reward Sparsity. Preference data often contains responses that, despite being preferred in pairwise comparisons, exhibit substantial variation in quality. As a result, the optimal responses—those associated with the highest reward value r_{max} —are sparse in the dataset. Since direct alignment algorithms do not account for these reward values, the trained model struggles to differentiate between responses of varying quality and fails to generalize effectively to the sparse optimal responses.

3.2 REWARD-CONDITIONED POLICIES LEARN FROM THE FULL SPECTRUM OF RESPONSE QUALITY

A straightforward way to address the limitations of direct alignment algorithms—specifically, their inability to account for the quality of responses—is to optimize a reward-conditioned policy. In this

¹For simplicity, we write $(x, y_w, y_l) \in \mathcal{D}_N$ as $y_w \succ y_l$.

162 163	response	y_1	y_2
164	r(x,y)	9	8
165	$\mathcal{D}_{N=1}$	$\{y_1\}$	$> y_2 \}$
166	$\pi^*(y \mid x) \qquad $	1	0
167	$\pi^*(u \mid x, a = 9) \mid$	1	0
168	$\pi^{*}(y \mid x, g = 8)$	0	1
169		-	

Table 1: High-quality rejected responses such as y_2 can be unnecessarily unlearned: $\pi^*(\cdot | x)$ deterministically generates y_1 . Rewardconditioned policies learn both responses and are easier to generalize to g = 10 with the extracted features from g = 8 and g = 9.

response	y_1	y_2	y_3
r(x,y)	9	1	0
$\mathcal{D}_{N=2}$	$\{y_1 >$	$> y_3, y_2 >$	$\rightarrow y_3$ }
$\pi^*(y \mid x) \qquad \mid$	1-a	a	0
$\pi^*(y \mid x, g = 9) \mid$	1	0	0
$\pi^*(y \mid x, g = 1) \mid$	0	1	0
$\pi^*(y \mid x, g = 0) \mid$	0	0	1

Table 2: Low-quality chosen responses such as y_2 can be learned: π^* indiscriminately generates y_1 and y_2 . Reward-conditioned policies distinguish the differences and learn the behaviors corresponding to different reward scores.

approach, the LLM policy is trained to generate responses corresponding to different reward values, enabling it to become aware of and adapt to these reward distinctions. By doing so, the LLM not only learns the patterns associated with the preferred responses but also retains the valuable information from the rejected ones, preventing the unlearning of high-quality rejected responses. For example, in Table 1, reward-conditioned policies learn to generate both y_1 and y_2 , instead of unlearning y_2 . This reward-based conditioning also enhances the model's ability to differentiate between responses of varying quality, even if both are preferred over a rejected alternative, as illustrated in Table 2. Besides, by extracting common patterns across responses with different quality levels, the LLM becomes more generalizable and is capable of generating the highest-quality responses with reward r_{max} (e.g., $r_{max} = 10$), which are often sparse in the training preference data.

4 Method

170

171

172

173

174

175 176 177

178

179

181

182

183

185

186 187 188

189

196

197

199

With the above motivation, we propose a data relabeling method that constructs a reward-augmented dataset by conditioning the preference pairs on the reward values given by the judge model r. Specifically, we define the goal-conditioned reward function $R(x, y, g) = -(g - r(x, y))^2$ as a function of the reward function r. The objective of the reward-conditioned policy $\pi(y \mid x, g)$ is thus to minimize the square difference between the goal reward g and the response reward r(x, y), which is equivalent to maximizing the goal-conditioned reward R(x, y, g), i.e.,

$$\min_{x \in \mathcal{D}_N, y \sim \pi(\cdot|x,g)} \left[(g - r(x,y))^2 \right] = \max_{x \in \mathcal{D}, y \sim \pi(\cdot|x,g)} \left[R(x,y,g) \right].$$
(4.1)

To optimize the RHS of Equation (4.1), we first observe that under the new goal-conditioned reward metric r, for each preference pair $x^i, y^i_w, r^i_w, y^i_l, r^i_l$ in \mathcal{D}_N , we have

$$\begin{split} R(x, y_w^i, g = r_w^i) &= 0 > R(x, y_l^i, g = r_w^i) = -(r_w^i - r_l^i)^2, \\ R(x, y_l^i, g = r_l^i) &= 0 > R(x, y_w^i, g = r_l^i) = -(r_w^i - r_l^i)^2. \end{split}$$

Thus, each pair can be relabeled to create two new preference pairs based on two distinct goals: when $g = r_w^i$, $y_w^i > y_l^i$; when $g = r_l^i$, $y_l^i > y_w^i$. Then any direct alignment algorithm can be applied to this new goal-conditioned preference dataset. Compared to fine-tuning on the original dataset \mathcal{D}_N , the model learns to capture not only desirable behaviors but also undesirable ones from the reward-augmented dataset. This approach helps identify patterns across high- and low-quality responses, enabling the LLMs to discern and learn from the entire spectrum of response quality and extrapolate to more optimal responses at inference time, by conditioning on higher reward goals.

We illustrate our method in Figure 1. For each preference pair with index i in \mathcal{D}_N , two goals are defined, corresponding to the reward values of the chosen response y_w^i and the rejected response y_l^i . Specifically, under the first goal $g = r_w^i$, the relabeled rewards are $R(x, y_w^i, g) = 0$ and $R(x, y_l^i, g) = -(r_w^i - r_l^i)^2$. The original ranking of responses remains the same, except that the LLM is preference optimized conditioned on $g = r_w^i$. Similarly, under the second goal $g = r_l^i$, the relabeled rewards are $R(x, y_l^i, g) = 0$ and $R(x, y_w^i, g) = -(r_w^i - r_l^i)^2$. Thus, the chosen and rejected



Figure 1: Construction of the reward-augmented preference dataset.

responses are reversed as y_l^i and y_w^i , respectively. By generating preference pairs conditioned on the goal reward for both the chosen and rejected responses, we obtain a reward-augmented dataset of size 2N. Finally, this new dataset can be used with any direct alignment algorithm, such as DPO.

In this work, we implement the reward-conditioned policy $\pi(y \mid x, g)$ as the LLM with a system prompt (or a prefix before the user prompt x if system prompts are not supported by the LLM) such as "generate responses of score g". At inference time, the LLM is conditioned on the optimal goal $g^* = r_{\text{max}}$ that is the highest possible reward value, e.g., $g^* = r_{\text{max}} = 10$, to generate the responses.

We provide the following theoretical guarantees for our method (see A.4 for a formal description).

Theorem 4.1 (Informal). Let $J(\pi) = \mathbb{E}_{x \sim d_0, y \sim \pi(\cdot | x, g^*)} [R(x, y, g^*)]$ be the performance measure, where R denotes the ground-truth goal-conditioned reward function and g^* denotes the optimal goal. Under mild assumptions, the policy $\hat{\pi}$ optimized from the reward-augmented DPO with a SFT regularizer satisfies that with probability at least $1 - \delta$,

$$J(\pi^{*}) - J(\widehat{\pi}) \leq \sqrt{\frac{1}{N}} \cdot \left\{ \frac{\sqrt{6}}{4} (1 + \exp(B))^{2} ((C_{\mu_{\mathcal{D}}}(\mathcal{R}; \pi^{*}, \pi_{\mathrm{sft}}))^{2} + 1) \iota + \mathbb{E}_{x \sim d_{0}} \left[\mathrm{KL}(\pi^{*}(\cdot|x, g^{*}) || \pi_{\mathrm{ref}}(\cdot|x, g^{*})) \right] \right\},$$
(4.2)

246 where $\pi^* = \operatorname{argmax}_{\pi} J(\pi)$ and $\iota = \sqrt{\log (N_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})/\delta)}$ with $\varepsilon = (6 \cdot (1 + e^B) \cdot N)^{-1}$. Here, 247 N denotes the number of preference pairs in \mathcal{D} , B denotes the upper bound of the reward models, 248 and the partial coverage coefficient $C_{\mu_{\overline{D}}}(\mathcal{R}; \pi^*, \pi_{sft})$ is defined in Assumption A.3.

The detailed proof is provided in A.5. The above theorem shows that our method attains global convergence to the optimal policy and the suboptimality decays at the order of $N^{-1/2}$ (N denotes the size of the reward-augmented preference dataset), which provides a theoretical justification for the strong empirical performance of the introduced reward-augmented DPO. Unlike prior works on goal-conditioned RL with supervised learning (Yang et al., 2022; Ghosh et al., 2019), which typically establish weaker results such as local performance improvements or the optimization of a lower bound on $J(\pi)$, our analysis guarantees global convergence to the optimal policy. This distinction underscores the significance of integrating DPO-like methods with goal-conditioned approaches.

5 RELATED WORK

225 226 227

228

229

230

231

232

233

234

235 236

237

238

239

249

250

251

252

253

254

255

256

257 258

259

260 **Preference Dataset Construction.** In order for the LLMs to follow instructions and better align 261 with human intents, it is common practice to build a preference dataset containing a set of prompts 262 and a pair of responses for each prompt, whose qualities are ranked by humans (Ouyang et al., 263 2022) or judge models (Bai et al., 2022b). A popular pipeline (Cui et al., 2023; Tunstall et al., 264 2023b; Wang et al., 2024c; Ivison et al., 2023; Zhu et al., 2023a) for constructing offline (i.e., fixed) 265 datasets involves sampling off-policy responses from various LLMs for each prompt in the hope 266 to increase the response diversity. The preference data can also be generated online (Guo et al., 267 2024) or iteratively (Bai et al., 2022a; Xu et al., 2023; Gulcehre et al., 2023; Hoang Tran, 2024; Xiong et al., 2023; Dong et al., 2024; Calandriello et al., 2024; Rosset et al., 2024) by sampling and 268 ranking on-policy responses from the training LLM. Recent works (Zhang et al., 2024; Cen et al., 269 2024; Xie et al., 2024) have also proposed systematically exploring the responses online and actively eliciting the preference. The proposed method in this paper is orthogonal to the construction ways
 of the preference data and can be applied to any dataset created either off-policy or on-policy.

273 **Preference Optimization.** Preference optimization methods generally follow two approaches. 274 The first involves learning a point-wise reward model, such as the Bradley-Terry model, and using 275 RL algorithms like PPO (Schulman et al., 2017; Zheng et al., 2023; Xu et al., 2024b) or REIN-FORCE (Williams, 1992; Li et al., 2023; Ahmadian et al., 2024), to maximize the KL-regularized 276 expected reward. The second approach is direct alignment (Rafailov et al., 2024b; Azar et al., 2023; 277 Zhao et al., 2023; Ethayarajh et al., 2024; Liu et al., 2024), which gets rid of a separate reward model 278 that is computationally costly to train. In this work, we mainly focus on the limitations of direct 279 alignment algorithms, particularly their unawareness of the quality aspects of responses. For PPO-280 style alignment algorithms that fit and maximize an explicit reward, preference data is only used to 281 learn the reward model, and policy training is performed in an online manner, where responses are 282 sampled from the LLM and their reward values directly play a role during the RL optimization. This 283 avoids drawbacks inherent to direct alignment methods, as detailed in Section 3.1.

284

285 **Conditional LLM Fine-Tuning.** Conditioning LLMs during training has proven effective for 286 aligning responses with specific human objectives. SteerLM (Dong et al., 2023b; Wang et al., 2023b) 287 extends SFT by conditioning the LLM on the multi-dimensional annotated attributes in data, such 288 as humor and toxicity, in order to steer model responses with user customizability. Directional Preference Alignment (DPA) (Wang et al., 2024a) proposed a variant of rejection sampling fine-tuning 289 (Yuan et al., 2023; Dong et al., 2023a) that conditions on the direction of the multi-objective re-290 ward, i.e., a user-dependent linear combination of the reward attributes (helpfulness and verbosity 291 in their experiments), that represents diverse preference objectives. These methods aim to train a 292 single LLM that can flexibly adjust to various user preference profiles. On the contrary, our method 293 targets the limitations of direct alignment algorithms by introducing reward-augmented relabeling. 294 This also differs from Conditioned-RLFT (Wang et al., 2023a), which leverages the data source in-295 formation by learning a class-conditioned policy with RL-free supervised learning. Reward-aware 296 Preference Optimization (RPO), introduced in Nemotron-4 (Adler et al., 2024), attempts to approx-297 imate the reward gap using the implicit reward and is motivated to resolve the unlearning issues of 298 DPO, which our work also addresses. However, we show that more limitations beyond unlearn-299 ing can be simply fixed with reward-conditioned LLMs and propose an easy-to-implement data relabeling method that integrates seamlessly with any direct alignment algorithm. Notably, Noise 300 Contrastive Alignment (Chen et al., 2024) and Unified Language Model Alignment (Cai et al., 2023) 301 introduce unified frameworks for alignment with binarized or reward datasets by leveraging (infor-302 mation) noise contrastive estimation and a hybrid of SFT with point-wise DPO, respectively. In 303 contrast, our work focuses on addressing the limitations of direct alignment algorithms with data 304 relabeling (on implicit-reward augmented binarized or reward datasets), and do not make algorithm 305 changes. We compare with all the aforementioned methods in our experiments. 306

6 EXPERIMENTS

6.1 REWARD-AUGMENTED DATA BOOSTS DPO PERFORMANCE

We begin by conducting experiments to demonstrate that applying the proposed method to fixed offline preference datasets leads to consistent performance improvements in DPO.

Setup. We adopt the UltraFeedback (Cui et al., 2023) preference dataset containing reward values scored by GPT-4 (LLM-as-Judge) that is ranged between 1 and 10 for each of the preference pairs. Our method constructs reward-augmented data by conditioning on these judge values. We fine-tune on various open-weight LLMs, including Mistral-7B-Instruct-v0.3, Qwen2-7B-Instruct, Llama-3.1-8B-Instruct, Gemma-2-9B-It, and SPPO (fine-tuned from Gemma2-9B-It). We use the DPO implementation in the Huggingface Alignment Handbook (Tunstall et al.). The hyperparameters and prompts that we use are listed in Appendix B.1.

321

307

308 309

310

Results. We first report the performance of the trained models on instruction-following benchmarks that use LLM as a judge, including AlpacaEval 2.0 (Dubois et al., 2024), MT-Bench (Zheng et al., 2024), and Arena-Hard-Auto (Li et al., 2024b). The results are shown in Figure 2. Across all instruction-following benchmarks, we observe that LLMs fine-tuned with DPO on the proposed reward-augmented data consistently outperform both their base models and those finetuned using DPO on the original UltraFeedback dataset by a considerable margin. Notably, direct alignment with the original preference data can sometimes degrade the performance of base models on specific benchmarks, such as Arena-Hard-Auto, which involves more complex reasoning tasks. In contrast, alignment using the reward-augmented data consistently yields superior results not only due to the improved style format gained from performing DPO on UltraFeedback.



(a) AlpacaEval 2.0 results. Left: Length-Controlled (LC) win rates. Right: Win rates.



Figure 2: Performance of the base models, the models trained with DPO on UltraFeedback, and the
models trained with DPO on reward-augmented ultrafeedback on AlpacaEval 2.0, MT-Bench, and
Arena-Hard-Auto benchmarks. The complete table is deferred to Appendix B.2.

357

368

331 332

333

334

335

336

337

338

339

341

342 343

345

347 348

349 350

351 352

353

Besides, we also evaluate the models on academic multi-choice QA benchmarks, including GSM8K
(Cobbe et al., 2021), GPQA (Rein et al., 2023), MUSR (Sprague et al., 2023), TruthfulQA (Lin
et al., 2021), BBH (Suzgun et al., 2022), and ARC Challenge (Clark et al., 2018). To better reflect
the capabilities of LLMs, we adopt various settings for these benchmarks, including zero-shot, fewshot, and few-shot Chain-of-Thought (CoT). The results are shown in Table 3.

It can be observed that performing DPO on the reward-augmented preference data leads to better
average academic scores for most families of models compared to models fine-tuned on the original UltraFeedback dataset and the base models. Besides, we didn't observe severe alignment tax
phenomenons (Askell et al., 2021; Noukhovitch et al., 2024; Li et al., 2024a) after DPO, and our
method is able to improve the base models on most of the benchmarks.

369 6.2 ABLATION STUDIES

Our Method Improves the Utility of Preference Data. We provide two pieces of evidence that we method son set more inice out of the

our method can get more juice out of the preference data compared to directly applying DPO. Firstly, we evaluate SPPO (Wu et al., 2024) fine-tuned with DPO on UltraFeedback (UF). The results are shown in Table 4. Since the SPPO model is already trained on UltraFeedback from Gemma-2-9B-It, an additional round of DPO training

	LC WR	WR	MT	Arena
SPPO	55.60	49.61	8.40	47.6
+DPO (UF)	52.75	40.58	8.41	40.4
+DPO (RA)	60.97	66.41	8.73	49.0

Table 4: SPPO can be improved with DPO by performing reward augmentation on the same data.

070		C C L LOTT	675 A 1	1 11 10 1	- 1110		1 2 9	
3/0	Model	GSM8K	GPQA	MUSR	TruthfulQA	BBH	ARC	Average
379	Widdel	(8-s CoT)	(0-s)	(0-s)	(0-s)	(3-s)	(25-s)	Average
380	Mistral-7B-Instruct-v0.3	52.39	30.62	47.35	59.71	46.64	58.53	49.21
381	+DPO (UltraFeedback)	53.22	28.94	47.35	64.74	47.46	60.32	50.34
202	+DPO (Reward-Augmented)	51.86	28.02	46.56	65.90	46.36	61.60	50.05
302	Qwen2-7B-Instruct	78.24	32.80	44.58	57.31	55.20	53.75	53.65
383	+DPO (UltraFeedback)	78.17	32.80	44.31	58.91	54.49	53.75	53.74
384	+DPO (Reward-Augmented)	81.05	32.97	45.77	57.99	54.94	54.52	54.54
385	Llama-3.1-8B-Instruct	76.72	33.89	39.95	54.00	50.74	55.38	51.78
386	+DPO (UltraFeedback)	78.47	33.72	43.39	56.61	51.31	57.51	53.50
387	+DPO (Reward-Augmented)	78.77	32.55	43.52	63.32	51.57	56.48	54.37
307	Gemma-2-9B-It	81.35	36.33	46.03	60.15	59.42	64.85	58.02
388	+DPO (UltraFeedback)	83.32	34.14	46.56	65.12	59.78	66.41	59.22
389	+DPO (Reward-Augmented)	83.62	35.74	48.15	65.27	59.82	65.87	59.75
390	SPPO	79.83	35.91	44.97	62.56	59.61	63.74	57.77
391	+DPO (UltraFeedback)	81.73	33.64	45.50	65.72	59.16	66.89	58.77
392	+DPO (Reward-Augmented)	80.67	36.16	48.68	67.39	58.88	65.53	59.55

Table 3: Performance comparison between the LLMs after DPO on UltraFeedback, on rewardaugmented UltraFeedback, and their base models on academic multi-choice QA benchmarks in standard zero-shot, few-shot, and CoT settings. Here, n-s refers to n-shot, the **bold** texts represent the best results in each family of models.

with the same data significantly degrades its performance. In contrast, performing DPO on Reward Augmented (RA) UltraFeedback results in substantial performance gains for SPPO, indicating that
 our method enhances the utility of the preference data.

The second evidence is that after DPO, the implicit reward can be used to relabel and augment the same preference data. Specifically, after training Qwen2-7B-Instruct with DPO on UltraFeedback, we leverage the resulting model π_{DPO} to calculate the implicit reward for each prompt x and response y, i.e., $\hat{r} = \beta (\log \pi_{\text{DPO}}(y \mid x) - \log \pi_{\text{Qwen}}(y \mid x))$. Then we perform DPO on Qwen2-7B-Instruct

Implicit-Reward-Augmented using the 406 (IRA) UltraFeedback. The results are shown 407 in Table 5. We observe that augmenting the 408 data with the implicit reward from the DPO 409 (UF) model leads to superior performance 410 even compared to augmenting the data with 411 reward scores from the LLM judge, i.e., 412 DPO (RA). This result highlights that DPO does not fully exploit the potential of the 413 data. Moreover, this ablation demonstrates 414 that our method is compatible with binarized 415 preference datasets that only contain chosen 416 and rejected response pairs, bypassing the 417 need for reward scores from judge models. 418

	LC WR	WR	MT	Arena
Qwen2-7B-It	20.93	18.22	7.90	24.3
+DPO (UF)	21.46	19.35	8.33	21.9
+DPO (RA)	31.17	27.58	8.47	30.1
+DPO (IRA)	32.61	29.15	8.49	28.3

Table 5: A second round of DPO on the rewardaugmented data, i.e., DPO (IRA), relabeled with the implicit reward from the DPO model at the first round, i.e., DPO (UF), significantly improves it. Our method helps get more juice out of the *binarized* (i.e., without judge model rewards) preference data.

419 420

394

395

396

397 398

Reward-Augmented Data is Superior Not Just Due to Its Increased Size. In this part, we

show that the success of our method is 421 not merely due to the increased size of 422 the training dataset. To illustrate this, we 423 perform DPO on the dataset where re-424 ward augmentation is applied to the first 425 half of the UltraFeedback data, which we 426 denote as DPO (Half RA). By doing so, 427 the reward-augmented data is of the same 428 size as the original dataset, but with only 429 half of the prompts and the corresponding responses being utilized. It can be ob-430 served from Table 6 that DPO (Half RA) 431 outperforms fine-tuning on the whole Ul-

	LC WR	WR	MT	Arena
Qwen2-7B-It	20.93	18.22	7.90	24.3
+DPO (UF)	21.46	19.35	8.33	21.9
+DPO (RA)	31.17	27.58	8.47	30.1
+DPO (Half RA)	29.56	28.30	8.33	26.9
Gemma-2-9B-It	49.20	37.58	8.54	42.8
+DPO (UF)	50.70	35.02	8.54	35.8
+DPO (RA)	59.27	54.56	8.59	43.9
+DPO (Half RA)	53.12	43.74	8.66	41.3

Table 6: DPO trained on only half of the data with reward augmentation outperforms the baseline. traFeedback (UF) by a large margin and achieves comparable performance to applying reward augmentation across the entire UF dataset, which is denoted as DPO (RA).

435 **Reward-Augmented Data Mitigates the Unlearning Issue.** 436 We first demonstrate that DPO suffers from the limitation of unnecessarily unlearning high-quality rejected responses, as 437 discussed in Section 3.1. Specifically, on the test set of Ultra-438 Feedback, we calculate the log probability of each rejected 439 response for the Qwen2-7B-Instruct model, its DPO (UF) 440 model, and our method DPO (RA). In Figure 3, we plot the 441 expected log probability for rejected responses with reward 442 scores ≥ 5 . We find that DPO substantially decreases the 443 probability of these high-quality rejected responses, confirm-444 ing that the unlearning issue arises in practice. In contrast, our 445 method alleviates this issue, although the probability is still 446 slightly lower than the base model, which is proven to be the



Figure 3: Our method helps mitigate the unlearning issue of DPO.

feature of DPO (Rafailov et al., 2024a; Zhang et al., 2024; Xu et al., 2024b).

449 **Impact of the Accuracy of AI Feedback.** We consider the 19.8k prompts from a 1/3 subset of UltraFeedback following the setup from Snorkel (Hoang Tran, 2024). Five on-policy responses are 450 first generated from Llama-3-8B-Instruct. An external reward model is followed to rank these re-451 sponses. We choose the best and worst responses as the chosen and rejected ones. DPO is then 452 performed on the resulting preference pairs and the reward-augmented pairs. To ablate how our 453 method will be impacted by the accuracy of AI feedback, we experiment with two reward mod-454 els as the ranker: PairRM (Jiang et al., 2023b) and ArmoRM (Wang et al., 2024b). PairRM is a 455 small-sized (0.4B) pairwise reward model, while ArmoRM is a 8B model that is state-of-the-art 456 on RewardBench (Lambert et al., 2024) and much stronger than PairRM. We implement a variant 457 (denoted as RA+) of the proposed reward augmentation method that only conditions on the goal 458 rewards of the chosen responses, not those of the rejected ones, leading to same-sized datasets. 459

	Llama-3-		PairRM (0.4B)		Armol	RM (8B)
	8B-Instruct	DPO (UF)	DPO (RA+)	DPO (RA)	DPO (UF)	DPO (RA+)
LC WR	22.92	41.76	44.72	48.20	42.32	48.73
WR	23.15	45.79	44.70	53.17	42.79	45.36

Table 7: Ablation on the impact of AI feedback quality on the AlpacaEval 2.0 benchmark.

The results in Table 7 demonstrate that training on augmented data conditioned on both chosen and rejected rewards is necessary for PairRM feedback, while relabeling with only the chosen rewards is sufficient to achieve strong performance for ArmoRM feedback. This aligns with our motivation outlined in Section 3.1: in noisy preference data, rejected responses may actually be of high quality, unlearning which can degrade performance. Similarly, low-quality chosen responses may also be reinforced. This issue does not arise with strong reward models that provide accurate preferences.

	SLiC-HF	ORPO	CPO	RRHF	KTO	IPO	RPO	R-DPO	SimPO	Ours
LC WR	26.9	28.5	28.9	31.3	33.1	35.6	40.8	41.1	44.7	48.2
WR	27.5	27.4	32.2	28.4	31.8	35.6	41.7	37.8	40.5	53.2

Table 8: Comparison between our method, i.e., Llama-3-8B-Instruct+DPO (RA) and baselines finetuned on the same model and on-policy data ranked by PairRM.

Moreover, in Table 8, we compare our method and various baselines under the same setting on the
AlpacaEval 2.0 benchmark, including SLiC-HF (Zhao et al., 2023), ORPO (Hong et al., 2024), CPO (Xu et al., 2024a), RRHF (Yuan et al., 2024), KTO (Ethayarajh et al., 2024), IPO (Azar et al., 2023),
R-DPO (Park et al., 2024), and SimPO (Meng et al., 2024), where the results are from Meng et al.
(2024), as well as the RPO (Adler et al., 2024) baseline that we implement. Our method outperforms the above algorithms by a considerable margin.

485

448

465 466

467

468

469

470

471

477

478

Conditioning on Multi-Attribute Rewards Enables SOTA Models. In previous parts, our

486 method is implemented by conditioning 487 on the scalar reward values given by 488 the judge models, either LLMs or re-489 ward models. We find that our approach 490 is generalizable to settings of multidimensional rewards that correspond to 491 different attributes, such as helpfulness 492 and truthfulness. Specifically, we follow 493

	LC Win Rate	Win Rate	Avg. Len.
Ours	56.57	52.19	1840
SimPO	53.70	47.50	1777
OpenChat	17.48	11.36	1362

Table 9: Our method trained with DPO achieves SOTA when conditioning on 5-dim rewards.

the setting from last part to construct the preference dataset by applying the ArmoRM reward model 494 on the on-policy responses generated by Llama-3-8B-Instruct. Since ArmoRM is a multi-objective 495 model that not only gives a scalar reward value but also predicts human-interpretable fine-grained 496 attributes, we first select 5 attributes (namely complexity, instruction following, honesty, helpful-497 ness, and intelligence depth) that have the highest average coefficients on the UltraFeedback data. 498 Then we relabel the data by conditioning on the 5-dim reward and follow the implementation of 499 using ArmoRM described in the last part. The resulting model achieves state-of-the-art within the 500 Llama-3-8B-Instruct model family, surpassing the strong baselines including SimPO (Meng et al., 2024) that is trained also on on-policy data ranked by ArmoRM, and OpenChat (Wang et al., 2023a) 501 fine-tuned with Conditioned-RLFT from the same Llama-3-8B-Instruct model. 502

Comparison with Conditional Fine-Tuning Baselines. We further compared with additional 504 conditional post-training baselines on the offline UltraFeedback dataset (i.e., without on-policy 505 data), including DPA (Wang et al., 2024a), SteerLM (Dong et al., 2023b), and (Info)NCA (Chen 506 et al., 2024). Since both baselines aim to optimize a user-controllable attribute-conditioned LLM 507 that is optimal under diverse preference profiles with different coefficients of the reward's attributes, 508 in Figure 4, we plot the win rate curves of these methods under varying preference profiles, such 509 as adjusting verbosity preferences as considered in Wang et al. (2024a). It can be observed that 510 fine-tuned from Zephyr-SFT, our method achieves the best AlpacaEval 2.0 win rate. In addition to the comparison with the implemented RPO (Adler et al., 2024) in Table 8, we also report the perfor-511 mance of RPO fine-tuned on additional models including Qwen2-7B-Instruct and Gemma2-9B-It. 512 As shown in Table 10, the implemented RPO is outperformed by our method across these models. 513



	LC Win Rate	Win Rate	Avg. Len.
Qwen+RPO	20.29	17.34	1704
Qwen+DPO (RA)	31.17	27.58	1789
Gemma+RPO	43.14	30.93	1413
Gemma+DPO (RA)	59.27	54.56	1872

Table 10: Comparison on AlpacaEval 2.0 between our method and RPO fine-tuned from the Qwen2-7B-Instruct and Gemma2-9B-It models. Our method consistently outperforms RPO across these fine-tuned models.

Figure 4: Comparison with DPA, SteerLM, and (Info)NCA.

7 CONCLUSION

514

515

516

517 518

519 520

521

522

523

524

526

In this paper, we first investigate the limitations of direct alignment algorithms, which arise from fo-527 cusing solely on relative preferences while neglecting the qualities of the responses and their gaps. 528 Specifically, since many rejected responses are only slightly worse than the chosen ones, striving to 529 maximize the reparameterized reward gap will cause overfitting and unnecessarily suppressing the 530 high-quality rejected response. Moreover, the directly aligned LLMs often struggle to differentiate 531 between responses of varying quality, indiscriminately learning the low-quality chosen responses 532 and failing to generalize effectively to more optimal responses that are sparse in the preference 533 data. To resolve the above limitations, we introduce a straightforward solution-learning reward-534 conditioned policies. By optimizing the LLM to generate responses conditioned on their qualities, it can better differentiate between quality levels and learn from the entire spectrum. Motivated by this, we propose a data relabeling method that constructs reward-augmented datasets by conditioning on 537 the quality of responses as the goal quality. In experiments, we fine-tune various LLMs by applying DPO on our reward-augmented data. The results demonstrate that our approach consistently de-538 livers significant performance improvements across various instruction-following benchmarks and increases the average accuracy on academic benchmarks.

540 REFERENCES

578

579

580

581

588

- Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H Anh, Pallab Bhattacharya, Annika Brundyn, Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, et al. Nemotron-4 340b technical report. *arXiv preprint arXiv:2406.11704*, 2024. 1, 2, 6, 9, 10
- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Ahmet Üstün, and
 Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human
 feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024. 6
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*, 2021. 7
- Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal
 Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human
 preferences. *arXiv preprint arXiv:2310.12036*, 2023. 1, 3, 6, 9
- 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*, 2022a. 5
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b. 2, 5
- 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. 16
- Tianchi Cai, Xierui Song, Jiyan Jiang, Fei Teng, Jinjie Gu, and Guannan Zhang. Ulma: Unified
 language model alignment with demonstration and point-wise human preference. *arXiv preprint arXiv:2312.02554*, 2023. 6
- Daniele Calandriello, Daniel Guo, Remi Munos, Mark Rowland, Yunhao Tang, Bernardo Avila
 Pires, Pierre Harvey Richemond, Charline Le Lan, Michal Valko, Tianqi Liu, et al. Human
 alignment of large language models through online preference optimisation. *arXiv preprint arXiv:2403.08635*, 2024. 5
- Shicong Cen, Jincheng Mei, Katayoon Goshvadi, Hanjun Dai, Tong Yang, Sherry Yang, Dale Schuurmans, Yuejie Chi, and Bo Dai. Value-incentivized preference optimization: A unified approach to online and offline rlhf. *arXiv preprint arXiv:2405.19320*, 2024. 5, 16, 17
- Huayu Chen, Guande He, Lifan Yuan, Ganqu Cui, Hang Su, and Jun Zhu. Noise contrastive alignment of language models with explicit rewards. *arXiv preprint arXiv:2402.05369*, 2024. 6, 10
 - Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018. 7
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 7
- 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. *arXiv* preprint arXiv:2310.01377, 2023. 2, 5, 6
- Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao,
 Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023a. 6
- Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen
 Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf. arXiv e-prints, pp. arXiv-2405, 2024. 2, 5

623

628

629

630

634

635

636

- 594 Yi Dong, Zhilin Wang, Makesh Narsimhan Sreedhar, Xianchao Wu, and Oleksii Kuchaiev. 595 Steerlm: Attribute conditioned sft as an (user-steerable) alternative to rlhf. arXiv preprint 596 arXiv:2310.05344, 2023b. 6, 10 597
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 598 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024. 2 600
- 601 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. arXiv preprint arXiv:2404.04475, 2024. 602 2,6 603
- 604 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model 605 alignment as prospect theoretic optimization. arXiv preprint arXiv:2402.01306, 2024. 1, 3, 6, 9 606
- Dibya Ghosh, Abhishek Gupta, Ashwin Reddy, Justin Fu, Coline Devin, Benjamin Eysenbach, 607 and Sergey Levine. Learning to reach goals via iterated supervised learning. arXiv preprint 608 arXiv:1912.06088, 2019. 5, 18 609
- 610 Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek 611 Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. Reinforced self-training 612 (rest) for language modeling. arXiv preprint arXiv:2308.08998, 2023. 5
- 613 Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre 614 Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from 615 online ai feedback. arXiv preprint arXiv:2402.04792, 2024. 5 616
- 617 Braden Hancock Hoang Tran, Chris Glaze. Snorkel-mistral-pairrm-dpo. 2024. URL https: //huggingface.co/snorkelai/Snorkel-Mistral-PairRM-DPO. 5,9 618
- 619 Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without 620 reference model. arXiv preprint arXiv:2403.07691, 2(4):5, 2024. 9 621
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep 622 Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. Camels in a changing climate: Enhancing lm adaptation with tulu 2. arXiv preprint arXiv:2311.10702, 2023. 5 624
- 625 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 626 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 627 Mistral 7b. arXiv preprint arXiv:2310.06825, 2023a. 2
 - Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. arXiv preprint arXiv:2306.02561, 2023b. 9
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, 631 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward 632 models for language modeling. arXiv preprint arXiv:2403.13787, 2024. 9 633
 - Shengzhi Li, Rongyu Lin, and Shichao Pei. Multi-modal preference alignment remedies regression of visual instruction tuning on language model. arXiv preprint arXiv:2402.10884, 2024a. 7
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gon-637 zalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and 638 benchbuilder pipeline. arXiv preprint arXiv:2406.11939, 2024b. 2, 6 639
- 640 Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. Remax: A 641 simple, effective, and efficient reinforcement learning method for aligning large language models. 642 In Forty-first International Conference on Machine Learning, 2023. 6
- 643 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human 644 falsehoods. arXiv preprint arXiv:2109.07958, 2021. 7 645
- Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and 646 Zhaoran Wang. Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adver-647 sarial regularizer. arXiv preprint arXiv:2405.16436, 2024. 6, 15, 16, 17, 18, 19, 23

648 649	I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 23
650	Yu Meng, Mengzhou Xia, and Danqi Chen. Simple preference optimization with a
651	reference-free reward. arxiv preprint arXiv:2403.14734, 2024. 9, 10
652	Eric Mitchell. A note on dpo with noisy preferences & relationship to ipo, 2023. 23
654	Michael Neulaboutch Commel Louis Floring Struck and Assess C. Commille Longuese model
655	Alignment with elastic reset Advances in Neural Information Processing Systems 36, 2024.
656	angument with elastic reset. <i>Navances in Ivearal Information Processing Systems</i> , 50, 2024.
657	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
658	Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-
659	10w instructions with numan feedback. Advances in neural information processing systems, 55: 27730–27744, 2022, 2, 5
660	21130 21144, 2022. 2, 5
661	Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality
662	in direct preference optimization. arXiv preprint arXiv:2403.19159, 2024. 9
663	Rafael Rafailov, Joev Hejna, Ryan Park, and Chelsea Finn. From r to a^* : Your language model is
665	secretly a q-function. arXiv preprint arXiv:2404.12358, 2024a. 9
666	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
667	Finn. Direct preference optimization: Your language model is secretly a reward model. Advances
668	in Neural Information Processing Systems, 36, 2024b. 1, 2, 3, 6, 19
669	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-
671	rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench-
672	mark. arXiv preprint arXiv:2311.12022, 2023. 7
673	Corby Rosset Ching-An Cheng Arindam Mitra Michael Santacroce Ahmed Awadallah and
674	Tengyang Xie. Direct nash optimization: Teaching language models to self-improve with general
675	preferences. arXiv preprint arXiv:2404.03715, 2024. 5
676	Lake Schulmen Eilin Welchi Desfulle Dharingh Ales Dedfend and Oles Klimen Dravinghesling
677	optimization algorithms. arXiv preprint arXiv:1707.06347, 2017. 6
678	
679 680	Zayne Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. Musr: Testing the limits of chain-of-thought with multistep soft reasoning. <i>arXiv preprint arXiv:2310.16049</i> , 2023. 7
681	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,
682	Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks
684	and whether chain-of-thought can solve them. arXiv preprint arXiv:2210.09261, 2022. 7
685	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-
686	patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma
687	2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> , 2024.
688	2
689	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Shengyi Huang, Kashif Rasul,
690	Alvaro Bartolome, Alexander M. Rush, and Thomas Wolf. The Alignment Handbook. URL
692	https://github.com/huggingface/alignment-handbook.6
693	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Shengyi Huang, Kashif Rasul,
694	Alexander M. Rush, and Thomas Wolf. The alignment handbook. https://github.com/
695	huggingface/alignment-handbook,2023a.23
696	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Raiani, Kashif Rasul, Younes Belkada
697	Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct
698	distillation of lm alignment. arXiv preprint arXiv:2310.16944, 2023b. 2, 3, 5
699	Guan Wang Sijia Chang Yianyuan Zhan Yiangang Li San Song and Yang Liu. Openshet: Ad
700 701	vancing open-source language models with mixed-quality data. <i>arXiv preprint arXiv:2309.11235</i> , 2023a. 6, 10

702 Haoxiang Wang, Yong Lin, Wei Xiong, Rui Yang, Shizhe Diao, Shuang Qiu, Han Zhao, and Tong 703 Zhang. Arithmetic control of llms for diverse user preferences: Directional preference alignment 704 with multi-objective rewards. arXiv preprint arXiv:2402.18571, 2024a. 6, 10 705 Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences 706 via multi-objective reward modeling and mixture-of-experts. arXiv preprint arXiv:2406.12845, 707 2024b. 9 708 709 Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David 710 Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring 711 the state of instruction tuning on open resources. Advances in Neural Information Processing 712 Systems, 36, 2024c. 5 713 Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, 714 Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, et al. Helpsteer: Multi-715 attribute helpfulness dataset for steerlm. arXiv preprint arXiv:2311.09528, 2023b. 6 716 717 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine learning, 8:229-256, 1992. 6 718 719 Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play 720 preference optimization for language model alignment. arXiv preprint arXiv:2405.00675, 2024. 721 2,7722 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent 723 pessimism for offline reinforcement learning. Advances in neural information processing systems, 724 34:6683-6694, 2021. 17 725 726 Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and 727 Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q*-approximation 728 for sample-efficient rlhf. arXiv preprint arXiv:2405.21046, 2024. 5 729 Wei Xiong, Hanze Dong, Chenlu Ye, Han Zhong, Nan Jiang, and Tong Zhang. Gibbs sam-730 pling from human feedback: A provable kl-constrained framework for rlhf. arXiv preprint 731 arXiv:2312.11456, 2023. 5, 16, 17 732 733 Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. 734 Iterative preference learning from human feedback: Bridging theory and practice for rlhf under 735 kl-constraint. In Forty-first International Conference on Machine Learning, 2024. 2 736 Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton 737 Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of llm 738 performance in machine translation. arXiv preprint arXiv:2401.08417, 2024a. 9 739 740 Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, and Jason Weston. Some things are more cringe than 741 others: Preference optimization with the pairwise cringe loss. arXiv preprint arXiv:2312.16682, 2023. 5 742 743 Shusheng Xu, Wei Fu, Jiaxuan Gao, Wenjie Ye, Weilin Liu, Zhiyu Mei, Guangju Wang, Chao Yu, 744 and Yi Wu. Is dpo superior to ppo for llm alignment? a comprehensive study. arXiv preprint 745 arXiv:2404.10719, 2024b. 6, 9 746 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, 747 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint 748 arXiv:2407.10671, 2024. 2 749 750 Rui Yang, Yiming Lu, Wenzhe Li, Hao Sun, Meng Fang, Yali Du, Xiu Li, Lei Han, and Chongjie 751 Zhang. Rethinking goal-conditioned supervised learning and its connection to offline rl. arXiv 752 preprint arXiv:2202.04478, 2022. 5, 18 753 Chenlu Ye, Wei Xiong, Yuheng Zhang, Nan Jiang, and Tong Zhang. A theoretical analysis of 754 nash learning from human feedback under general kl-regularized preference. arXiv preprint 755 arXiv:2402.07314, 2024. 17

756 757 758	Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback. <i>Advances in Neural Information Processing Systems</i> , 36, 2024. 9
759 760 761 762	Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou, and Jingren Zhou. Scaling relationship on learning mathematical reasoning with large language models. <i>arXiv preprint arXiv:2308.01825</i> , 2023. 6
763 764 765	Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable offline preference-based reinforcement learning. In <i>The Twelfth International Conference on Learning Representations</i> , 2023. 15, 17
766 767 768	Shenao Zhang, Donghan Yu, Hiteshi Sharma, Ziyi Yang, Shuohang Wang, Hany Hassan, and Zhao- ran Wang. Self-exploring language models: Active preference elicitation for online alignment. <i>arXiv preprint arXiv:2405.19332</i> , 2024. 3, 5, 9
769 770 771	Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. Slic-hf: Sequence likelihood calibration with human feedback. <i>arXiv preprint arXiv:2305.10425</i> , 2023. 1, 3, 6, 9
772 773 774 775	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in Neural Information Processing Systems</i> , 36, 2024. 2, 6
776 777 778	Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, et al. Secrets of rlhf in large language models part i: Ppo. arXiv preprint arXiv:2307.04964, 2023. 6
779 780	Ding-Xuan Zhou. The covering number in learning theory. Journal of Complexity, 18(3):739–767,
781	2002. 15, 21
782 783	Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. Starling-7b: Improving llm helpfulness and harmlessness with rlaif, November 2023a. 3, 5
784 785	Banghua Zhu, Jiantao Jiao, and Michael I Jordan. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. arXiv preprint arXiv:2301.11270, 2023b. 17
786	
787 788	A THEORY
789 790	In this section, we present the theoretical analysis for our proposed method.
791 792	A.1 CONCEPTS
793	We provide some useful concepts for the simplicity of later discussions.
794 795	• Hellinger distance $D_{\text{Hellinger}}(p q)$ between two probability density functions p and q defined on \mathcal{X} is defined as
796 797	$D_{T} = \frac{1}{2} \int \left(\sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx$
797	$D_{\text{Hellinger}}(p q) = \frac{1}{2} \int_{x \in \mathcal{X}} \left(\sqrt{p(x)} - \sqrt{q(x)} \right) \mathrm{d}x.$
799	• Total variation (TV) distance $D_{TV}(p q)$ between two probability density functions p and
800	q defined on \mathcal{X} is defined as
801	$D_{-}(\mathbf{r} \ \mathbf{r}) = \frac{1}{2} \int \mathbf{r}(\mathbf{r}) - \mathbf{r}(\mathbf{r}) d\mathbf{r}$
802	$D_{\mathrm{TV}}(p q) = \frac{1}{2} \int_{x \in \mathcal{X}} p(x) - q(x) \mathrm{d}x.$
803	• Kullback–Leibler (KL) divergence $KL(p a)$ between two probability density functions p
804	and q defined on \mathcal{X} is defined as
805	VI (all x) $\int \int p(x) dx$
807	$KL(p q) = \int_{x \in \mathcal{X}} \log\left(\frac{1}{q(x)}\right) p(x) \mathrm{d}x.$
808	• We denote $\mathcal{N}_{\epsilon}(\mathcal{F} \parallel \cdot \parallel_{\infty})$ as the ϵ -covering number (Zhou, 2002) for function class \mathcal{F} under
809	the infinity norm $\ \cdot\ _{\infty}$. Widely used in the theoretical analysis (Liu et al., 2024; Zhan et al.,

the infinity norm $\|\cdot\|_{\infty}^{\infty}$. Widely used in the theoretical analysis (Liu et al., 2024; Zhan et al., 2023), the ϵ -covering number characterizes the complexity of the function class \mathcal{F} .

810 A.2 THEORETICAL FORMULATION

Goal-conditioned preference model. Consider a language model $\pi \in \Delta_{\mathcal{Y}}^{\mathcal{X}}$ that takes the prompt $x \in \mathcal{X}$ as input and outputs the response $y \in \mathcal{Y}$, where \mathcal{X} and \mathcal{Y} are spaces of prompts and responses, respectively. Given the prompt $x \in \mathcal{X}$, a discrete probability distribution $\pi(\cdot \mid x) \in \Delta_{\mathcal{Y}}$ is generated, where $\Delta_{\mathcal{Y}}$ is the set of discrete distributions over \mathcal{Y} . We define the goal-conditioned reward function class as $\mathcal{R} \subset \{R(x, y, g) : \mathcal{X} \times \mathcal{Y} \times \mathcal{G} \mapsto \mathbb{R}\}$, where \mathcal{G} is the goal space. The goal-conditioned Bradley-Terry model (Bradley & Terry, 1952) for annotations is described as

$$\mathbb{P}_{R}(y_{1} \succ y_{0} | x, y_{1}, y_{0}, g) = \frac{\exp(R(x, y_{1}, g))}{\exp(R(x, y_{1}, g)) + \exp(R(x, y_{0}, g))} = \sigma(R(x, y_{1}, g) - R(x, y_{0}, g)),$$
(A.1)

where $\sigma(z) = 1/(1 + \exp(-z))$ is the sigmoid function. For notational simplicity, we also denote that the reward is parameterized by $\theta \in \Theta$. We denote the corresponding negative log-likelihood function for r for a reward-augmented preference dataset $\overline{\mathcal{D}} = \{(x^i, y^i_w, y^i_l, g^i)\}_{i=1}^N$ as

$$\mathcal{L}(R,\bar{\mathcal{D}}) = -\mathbb{E}_{(x,y_w,y_l,g)\sim\bar{\mathcal{D}}} \left[\log \sigma \left(R(x,y_w,g) - R(x,y_l,g) \right) \right], \tag{A.2}$$

where y_w^i is preferred to y_l^i by the annotation given the prompt x^i and goal g^i for any $i \in [N]$. For notational simplicity, we denote the DPO loss by

$$\mathcal{L}_{\text{DPO}}(\pi,\bar{\mathcal{D}}) = -\mathbb{E}_{(x,y_w,y_l,g)\sim\bar{\mathcal{D}}} \bigg[\log \sigma \bigg(\beta \log \frac{\pi(y_w \mid x,g)}{\pi_{\text{ref}}(y_w \mid x,g)} - \beta \log \frac{\pi(y_l \mid x,g)}{\pi_{\text{ref}}(y_l \mid x,g)} \bigg) \bigg].$$
(A.3)

Performance metric. For the notational simplicity in the theoretical analysis, we denote by R^* the ground-truth goal-conditioned reward function. The alignment target is to maximize the expected true reward R^* conditioned on the optimal goal $g^* \in \mathcal{G}$. Thus, we define the value function of any policy π as

$$J(\pi) = \mathbb{E}_{x \sim d_0, y \sim \pi(\cdot|x, g^\star)} \left[R^\star(x, y, g^\star) \right]. \tag{A.4}$$

Here we allow the prompt distribution $d_0(\cdot)$ to be different from that of the offline dataset distribution $\mu_{\overline{D}}(\cdot)$, but is assumed to be known. In the meanwhile, we consider the policies that share the same support as the reference policy π_{ref} (Xiong et al., 2023), that is, we take a policy class Π as

$$\Pi = \left\{ \pi : \mathcal{X} \times \mathcal{G} \mapsto \Delta(\mathcal{A}) \, \middle| \, \operatorname{Supp}(\pi(\cdot|x,g)) \subseteq \operatorname{Supp}(\pi_{\operatorname{ref}}(\cdot|x,g)), \, \forall (x,g) \in \mathcal{X} \times \mathcal{G} \right\}.$$
(A.5)

The performance gap of a learned policy $\hat{\pi} \in \Pi$ w.r.t. any given optimal policy π^* is measured as

$$\operatorname{Gap}^{\pi^{\star}}(\widehat{\pi}) = J(\pi^{\star}) - J(\widehat{\pi}), \text{ given any optimal policy } \pi^{\star} \in \Pi,$$
(A.6)

One popular choice to define the optimal policy is to maximize the KL-regularized reward, i.e.,

$$\pi^{\star} = \operatorname*{argmax}_{\pi \in \Pi} \left[R^{\star}(x, y, g^{\star}) - \beta_0 \mathrm{KL} \left(\pi(\cdot \mid x, g^{\star}) \| \pi_{\mathrm{ref}}(\cdot \mid x, g^{\star}) \right) \right]$$
(A.7)

for a fixed $\beta_0 > 0$.

Theoretical version of the reward-augmented DPO. We formulate the theoretical version of the reward-augmented DPO in Algorithm 8, where we add a SFT regularizer on the empirical objective to handle the issue distribution shift and analyze the bound on the suboptimality (Liu et al., 2024; Cen et al., 2024). One simple choice to define SFT policy π_{SFT} is to utilize the chosen labels in the original preference dataset D, that is,

$$\pi_{\text{sft}} = \operatorname*{argmax}_{\pi \in \Pi} \mathbb{E}_{(x, y_w) \sim \mathcal{D}}[\log \pi(y_w \mid x, g^*)].$$
(A.8)

861 In practice, the goal relabeling distribution $g \sim p_{\mathcal{G}}(\cdot | x, y)$ is set to be a deterministic selection 862 of the annotated reward of the chosen response, i.e., $g = r_{\text{RM}}(x, y)$ for any $i \in [N]$ and a given 863 reward model r_{RM} . We also remark that the size of the reward-augmented preference dataset $\overline{\mathcal{D}}$ is $N = 2N_0$, where N_0 denotes the size of the original preference dataset \mathcal{D} .

864

865

866

877 878

879 880

882

883 884

885

Algorithm 1 Theoretical Version of the Reward-Augmented DPO

1: **Input**: Preference dataset $\mathcal{D} = \{(x^i, y^i_w, y^i_l)\}_{i=1}^{N_0}$, parameters $\beta, \eta > 0$, reference policy π_{ref} , SFT policy π_{sft} for the regularizer, and goal labeling_distribution $p_{\mathcal{G}}$.

2: Initialize the reward-augmented preference dataset $\overline{D} = \emptyset$.

3: for $i = 1, ..., N_0$ do

4: Sample goal g_w^i from $p_{\mathcal{G}}(\cdot | x^i, y_w^i)$ and update the reward-augmented preference dataset as $\bar{\mathcal{D}} \leftarrow \bar{\mathcal{D}} \cup \{(x^i, y_w^i, y_l^i, g_w^i)\}.$

5: Sample goal g_l^i from $p_{\mathcal{G}}(\cdot | x^i, y_l^i)$ and update the reward-augmented preference dataset as $\overline{\mathcal{D}} \leftarrow \overline{\mathcal{D}} \cup \{(x^i, y_l^i, y_w^i, g_l^i)\}.$

6: end for

7: Solve policy $\pi_{\widehat{\alpha}}$ by optimizing the following objective

$$\min_{\theta \in \Theta} \left\{ \mathbb{E}_{x \sim d_0, y_0 \sim \pi_{\text{sft}}(\cdot | x, g^*)} \left[-\eta \beta \cdot \log(\pi_\theta(y_0 | x, g^*)) \right] + \mathcal{L}_{\text{DPO}}(\pi_\theta, \bar{\mathcal{D}}) \right\}$$
(A.9)

8: **Output**: Policy $\widehat{\pi} = \pi_{\widehat{\theta}}$.

A.3 ASSUMPTIONS FOR THEORETICAL ANALYSIS

Similar to the theoretical analyses on offline RLHF (Liu et al., 2024; Cen et al., 2024), we provide the following assumptions.

886 887 888 888 889 Assumption A.1 (True reward model). We assume that the true goal-conditioned reward model $R^* \in \mathcal{R}$ for, and for any $R \in \mathcal{R}$ and $(x, y, g) \in \mathcal{X} \times \mathcal{A} \times \mathcal{G}$, it holds that $R(x, y, g) \in [-B/2, B/2]$ for a positive constant B > 0.

Assumption A.1 is standard in sample complexity analysis (Zhu et al., 2023b; Zhan et al., 2023; Ye et al., 2024) in RLHF.

Assumption A.2 (Regularity). We assume that the reward model class \mathcal{R} , prompt space \mathcal{X} , and goal space \mathcal{G} are convex and compact.

Assumption A.2 plays a role in establishing the equivalence between maximin and minimax optimizations. This assumption is naturally satisfied when considering a linear reward function (Zhu et al., 2023b; Xiong et al., 2023; Cen et al., 2024) of the form $R_{\theta}(x, y, g) = \varphi(x, y, g)^{\top} \theta$, where φ represents a known feature map. More broadly, the assumption is also met by the class of Lipschitz continuous reward models.

Assumption A.3 (Partial coverage coefficient). Given the optimal policy $\pi^* \in \Pi$, the coverage coefficient of the population distribution $\mu_{\bar{D}}$ of the reward-augmented preference dataset \bar{D} w.r.t. reward model class \mathcal{R} , optimal policy π^* , and the SFT policy π_{sft} , denoted by $C_{\mu_{\bar{D}}}(\mathcal{R}; \pi^*, \pi_{\text{sft}})$, is defined as

$$\sup_{R \in \mathcal{R}} \frac{\mathbb{E}_{x \sim d_0, y_1 \sim \pi^*(\cdot | x, g^*), y_0 \sim \pi_{\text{sft}}(\cdot | x, g^*)} \left[\left(R^*(x, y_1, g^*) - R^*(x, y_0, g^*) - \left(R(x, y_1, g^*) - R(x, y_0, g^*) \right) \right]}{\sqrt{\mathbb{E}_{(x, y_w, y_l, g) \sim \mu_{\overline{D}}}} \left[\left| \left(R(x, y_w, g) - R(x, y_l, g) \right) - \left(R(x, y_w, g) - R(x, y_l, g) \right) \right|^2 \right]}$$
(A.10)

908 909 910

911

905 906 907

We assume that $C_{\mu_{\overline{D}}}(\mathcal{R}; \pi^*, \pi_{\text{sft}}) < +\infty$ for the given optimal policy $\pi^* \in \Pi$.

Assumption A.3 characterizes how well the dataset \overline{D} covers the optimal policy π^* and the SFT policy π_{sft} given the optimal goal g^* , instead of covering all the policies in the policy class. That is the reason why we call this assumption "partial coverage". Different variants of partial coverage assumptions are posed in previous literature (Liu et al., 2024; Cen et al., 2024; Zhan et al., 2023; Xie et al., 2021) that study offline RLHF and RL to characterize the distribution shift between the optimal policy and the offline dataset distribution. We remark that the quantity $C_{\mu_{\overline{D}}}(\mathcal{R}; \pi^*, \pi_{\text{sft}})$ is upper bounded by the density ratio $||d_0(\cdot) \otimes \pi^*(\cdot|\cdot, g^*) \otimes \pi_{\text{sft}}(\cdot|\cdot, g^*)/\mu_{\overline{D}}(\cdot, \cdot, \cdot, g^*)||_{\infty}$.

918 A.4 THEORETICAL RESULTS

920 Under assumptions introduced before, we are ready to give the theoretical result for Algorithm 1 in921 the following theorem.

Theorem A.4 (Suboptimality of Algorithm 8). Taking the policy class Π as (A.5), supposing that Assumptions A.1, A.2, and A.3 hold, and assuming that the reward model class \mathcal{R} has a finite ε -covering number under $\|\cdot\|_{\infty}$ -norm $\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty}) < +\infty$ with $\varepsilon = (6 \cdot (1 + e^B) \cdot N)^{-1}$. Setting

$$\eta = (1 + \exp(B))^{-2} \cdot \sqrt{24 \log\left(\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})/\delta\right)/N}, \quad \beta = 1/\sqrt{N}$$

in Algorithm 1. Then the output policy $\hat{\pi}$ of Algorithm 1 satisfies that with probability at least $1 - \delta$,

$$\operatorname{Gap}^{\pi^{\star}}(\widehat{\pi}) \leq \sqrt{\frac{1}{N}} \cdot \left\{ \frac{\sqrt{6}}{4} \left(1 + \exp(B) \right)^{2} \left(\left(C_{\mu_{\widehat{D}}}(\mathcal{R}; \pi^{\star}, \pi_{\mathrm{sft}}) \right)^{2} + 1 \right) \iota \right. \\ \left. + \mathbb{E}_{x \sim d_{0}} \left[\operatorname{KL}\left(\pi^{\star}(\cdot | x, g^{\star}) \| \pi_{\mathrm{ref}}(\cdot | x, g^{\star}) \right) \right] \right\},$$
(A.11)

where $\iota = \sqrt{\log (\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})/\delta)}$ with $\varepsilon = (6 \cdot (1 + e^B) \cdot N)^{-1}$. Here, N denotes the number of preference pairs in \mathcal{D} , R denotes the upper bound of the reward models, and the partial coverage coefficient $C_{\mu_{\overline{\mathcal{D}}}}(\mathcal{R}; \pi^*, \pi_{\text{sft}})$ is defined in Assumption A.3.

The detailed proof is provided in Appendix A.5. Theorem A.4 shows that our proposed reward-augmented DPO (Algorithm 1) can attain global convergence to the optimal policy and the sub-optimality decays at the order of $N^{-1/2}$ (N denotes the size of the reward-augmented preference dataset). Theorem A.4 provides a theoretical justification for the strong empirical performance of the reward-augmented DPO introduced in this paper. Unlike prior works on goal-conditioned rein-forcement learning with supervised learning (Yang et al., 2022; Ghosh et al., 2019), which typically establish weaker results such as local performance improvements or the optimization of a lower bound on $J(\pi)$, our analysis guarantees global convergence to the optimal policy. This distinction underscores the significance of integrating DPO-like methods with goal-conditioned approaches.

948 A.5 PROOF OF THEOREM A.4

Bridge Algorithm 1 to the maximin optimization. Motivated by Liu et al. (2024), we transform
the optimization objective in Algorithm 1 to a minimax optimization objective, and then to a maximum optimization objective, where the maximum optimization objective can be analyzed with tools in RL analysis.

Define the objective function $\phi(\pi, R)$ as

$$\phi(\pi, R) = \eta \cdot \mathbb{E}_{\substack{x \sim d_0, y_1 \sim \pi(\cdot | x, g^*) \\ y_0 \sim \pi_{\text{sft}}(\cdot | x, g^*)}} \Big[R(x, y_1, g^*) - R(x, y_0, g^*) \\ - \beta \cdot D_{\text{KL}} \big(\pi(\cdot | x, g^*) \| \pi_{\text{ref}}(\cdot | x, g^*) \big) \Big] + \mathcal{L}(R, \bar{\mathcal{D}}).$$
(A.12)

First, we prove that the derived policy $\hat{\pi}$ from Algorithm 1 satisfies

$$\widehat{\pi} \in \operatorname*{argmax}_{\pi \in \Pi} \phi(\widehat{R}, \pi), \quad \text{where} \quad \widehat{R} \in \operatorname*{argmin}_{R \in \mathcal{R}} \max_{\pi \in \Pi} \phi(\pi, R). \tag{A.13}$$

By the definition of the optimization objective $\phi(\pi, R)$ in (A.12), we have

$$\min_{R\in\mathcal{R}}\max_{\pi\in\Pi}\phi(\pi,R) = \min_{R\in\mathcal{R}}\left\{\eta\cdot\max_{\pi\in\Pi}\left\{\mathbb{E}_{x\sim d_0,y_1\sim\pi(\cdot|x,g^\star)}\left[R(x,y_1,g^\star) - \beta\cdot\operatorname{KL}(\pi(\cdot|x,g^\star)\|\pi_{\operatorname{ref}}(\cdot|x,g^\star))\right]\right\} - \eta\cdot\mathbb{E}_{x\sim d_0,y_0\sim\pi_{\operatorname{sh}}(\cdot|x,g^\star)}\left[R(x,y_0,g^\star)\right] + \mathcal{L}(R,\bar{\mathcal{D}})\right\}.$$
(A.14)

Then, we apply the following lemma to solve the inner maximization problem in (A.14).

Lemma A.5 (Oracle optimal KL-regularized policy). Given any reward model $R \in \mathcal{R}$, the optimal policy π_R to the maximization problem

$$\max_{\pi \in \Pi} \left\{ \mathbb{E}_{x \sim d_0, y \sim \pi(\cdot | x, g^\star)} \Big[R(x, y, g^\star) - \beta \cdot \mathrm{KL}\big(\pi(\cdot | x, g^\star) \| \pi_{\mathrm{ref}}(\cdot | x, g^\star)\big) \Big] \right\}.$$
(A.15)

is given by

$$\pi_R(\cdot|x,g) = \frac{1}{Z_R(x,g)} \cdot \pi^{\text{ref}}(\cdot|x,g) \cdot \exp\left(\beta^{-1}R(x,\cdot,g)\right), \qquad (A.16)$$
$$Z_R(x,g) = \int_{y \in \mathcal{Y}} \exp\left(\beta^{-1}R(x,y,g)\right) \mathrm{d}\pi_{\text{ref}}(y|x,g),$$

and correspondingly the optimal value of (A.15) is given by $(A.15) = \mathbb{E}_{x \sim d_0}[\beta \cdot \log(Z_R(x, g^*))].$

Proof of Lemma A.5. See the proof in Lemma 4.2 of Liu et al. (2024).

By Lemma A.5 and (A.14), we have

$$\min_{R \in \mathcal{R}} \max_{\pi \in \Pi} \phi(\pi, R) = \min_{R \in \mathcal{R}} \left\{ \beta \eta \cdot \log(Z_R(x, g^*)) - \eta \cdot \mathbb{E}_{x \sim d_0, y_0 \sim \pi_{\text{sft}}(\cdot | x, g^*)} \left[R(x, y_0, g^*) \right] + \mathcal{L}(R, \bar{\mathcal{D}}) \right\}.$$
(A.17)

From Lemma A.5, we know that given any reward model $R \in \mathcal{R}$, we can reparameterize it via its corresponding optimal goal-conditioned KL-regularized policy π_R (Rafailov et al., 2024b), that is,

$$R(x, \cdot, g) = \beta \cdot \log\left(\frac{\pi_R(\cdot|x, g)}{\pi^{\text{ref}}(\cdot|x, g)}\right) + \beta \cdot \log(Z_R(x, g)).$$
(A.18)

Plugging (A.18) into (A.19), we show that the optimization problem in Algorithm 1 relates to the minimax optimization problem on $\phi(\pi, R)$:

$$\min_{R \in \mathcal{R}} \max_{\pi \in \Pi} \phi(\pi, R) = \min_{R \in \mathcal{R}} \left\{ \eta \beta \cdot \mathbb{E}_{x \sim d_0, y_0 \sim \pi_{\text{sft}}(\cdot | x, g^*)} \left[\log \left(\frac{\pi_R(y_0 \mid x, g^*)}{\pi_{\text{ref}}(y_0 \mid x, g^*)} \right) \right] + \mathcal{L}_{\text{DPO}}(\pi_R, \bar{\mathcal{D}}) \right\}$$
$$= \min_{R \in \mathcal{R}} \left\{ \eta \beta \cdot \mathbb{E}_{x \sim d_0, y_0 \sim \pi_{\text{sft}}(\cdot | x, g^*)} \left[\log \left(\pi_R(y_0 \mid x, g^*) \right) \right] + \mathcal{L}_{\text{DPO}}(\pi_R, \bar{\mathcal{D}}) \right\}.$$
(A.19)

where the first equality uses the definition of DPO loss \mathcal{L}_{DPO} in (A.3). Since we know that $\hat{\pi} \in \max_{\pi \in \Pi} \phi(\hat{R}, \pi)$ and \hat{r} solves the minimization problem in (A.19), we know that $\hat{\pi} = \pi_{\hat{R}}$ by Lemma A.5.

¹⁰⁰⁹ Next, we show that the minimization problem $\phi(\pi, R)$ can be equivalently transformed into a maximization problem. Specifically, we will prove that the output policy $\hat{\pi}$ for the Algorithm 1 satisfies

$$\widehat{\pi} \in \operatorname*{argmax}_{\pi \in \Pi} \min_{R \in \mathcal{R}} \phi(\pi, R), \tag{A.20}$$

1014 which is implied by the following theorem.

Theorem A.6. For the policy class Π defined in (A.5) and the reward model class \mathcal{R} satisfying Assumption A.2, consider the following policy defined as

$$\pi_{\widehat{R}} \in \operatorname*{argmax}_{\pi \in \Pi} \phi(\widehat{R}, \pi), \quad \text{where} \quad \widehat{R} \in \operatorname*{argmin}_{R \in \mathcal{R}} \max_{\pi \in \Pi} \phi(\pi, R).$$
(A.21)

1020 Then the policy $\pi_{\hat{R}}$ also solves the following maximin optimization:

$$\pi_{\widehat{R}} \in \operatorname*{argmax}_{\pi \in \Pi} \min_{R \in \mathcal{R}} \phi(\pi, R).$$
(A.22)

Proof. Under Assumption A.1, we know that $\phi(\pi, R)$ is convex for $R \in \mathcal{R}$ and strongly concave for $\pi \in \Pi$. Applying Theorem 5.6 in Liu et al. (2024), we prove Theorem A.6.

1026 1027 1028 Suboptimality Decomposition. By the definitions of the optimization objective $\phi(\pi, R)$ in (A.12) and the suboptimality gap of $\hat{\pi}$ w.r.t. π^* in (A.6), we decompose the gap as

where the inequality follows the fact that $\hat{\pi}$ solves the maxmin optimization problem in (A.20).

 $\operatorname{Gap}_{\beta}^{\pi^{*}}(\widehat{\pi}) = \operatorname{Term}(A) + \operatorname{Term}(B) + \operatorname{Term}(C)$

 $\mathcal{L}(R^{\star}, \bar{\mathcal{D}}) - \mathcal{L}(R, \bar{\mathcal{D}})$

Analysis of Term (B) in (A.23). Note that

1082 1083

1084 1085

$$\begin{aligned} \text{Term (B)} \\ &= \eta^{-1} \cdot \min_{R \in \mathcal{R}} \phi(\hat{\pi}, R) \\ &\quad - \mathbb{E}_{x \sim d_0, y_1 \sim \hat{\pi}(\cdot | x, g^*), y_0 \sim \pi_{\text{sft}(\cdot | x, g^*)}} \Big[R^*(x, y_1, g^*) - R^*(x, y_0, g^*) - \beta \cdot \text{KL}(\hat{\pi}(\cdot | x, g^*) \| \pi_{\text{ref}}(\cdot | x, g^*)) \Big] \\ &\leq \mathbb{E}_{x \sim d_0, y_1 \sim \hat{\pi}(\cdot | x, g^*), y_0 \sim \pi_{\text{sft}}(\cdot | x, g^*)} \Big[R^*(x, y_1, g^*) - R^*(x, y_0, g^*) - \beta \cdot \text{KL}(\hat{\pi}(\cdot | x, g^*) \| \pi_{\text{ref}}(\cdot | x, g^*)) \Big] \\ &\quad + \eta^{-1} \cdot \mathcal{L}(R^*, \bar{\mathcal{D}}) \\ &\quad - \mathbb{E}_{x \sim d_0, y_1 \sim \hat{\pi}(\cdot | x, g^*), y_0 \sim \pi_{\text{sft}(\cdot | x, g^*)}} \Big[R^*(x, y_1, g^*) - R^*(x, y_0, g^*) - \beta \cdot \text{KL}(\hat{\pi}(\cdot | x, g^*) \| \pi_{\text{ref}}(\cdot | x, g^*)) \Big] \\ &= \eta^{-1} \cdot \mathcal{L}(R^*, \bar{\mathcal{D}}), \end{aligned}$$
(A.28)

where the inequality uses the fact that $R^* \in \mathcal{R}$ by Assumption A.1 and the definition of the optimization objective in (A.12).

Concluding the remaining proof. Combining (A.23), (A.27), and (A.28), we have

1100 1101

1102

$$\leq \max_{R \in \mathcal{R}} \left\{ \mathbb{E}_{\substack{x \sim d_{0}, y_{1} \sim \pi^{*}(\cdot | x, g^{*}), \\ y_{0} \sim \pi_{\text{sff}}(\cdot | x, g^{*})}} \left[\left(R^{*}(x, y_{1}, g^{*}) - R^{*}(x, y_{0}, g^{*}) \right) - \left(R(x, y_{1}, g^{*}) - R(x, y_{0}, g^{*}) \right) \right] \\ + \eta^{-1} \cdot \left(\mathcal{L}(R^{*}, \bar{\mathcal{D}}) - \mathcal{L}(R, \bar{\mathcal{D}}) \right) \right\} \\ + \beta \cdot \mathbb{E}_{x \sim d_{0}} \left[\operatorname{KL}(\pi^{*}(\cdot | x, g^{*}) \| \pi_{\text{ref}}(\cdot | x, g^{*})) - \operatorname{KL}(\widehat{\pi}(\cdot | x, g^{*}) \| \pi_{\text{ref}}(\cdot | x, g^{*})) \right].$$
(A.29)

1103 1104 1105

1106 1107 Next, we upper bound the right-hand side of (A.29) by relating the negative log-likelihood loss dif-1108 ference term to the reward difference term. Recall the definition of the goal-conditioned preference 1109 model \mathbb{P}_R in (A.1). Applying Lemma A.7 to give an upper bound of the difference of the negative 1109 log-likelihood loss and setting $\varepsilon = (6 \cdot (1 + e^B) \cdot N)^{-1}$, it holds with probability at least $1 - \delta$ and 1110 for any reward model $R \in \mathcal{R}$ that

1116

1129

 $\leq -2 \cdot \mathbb{E}_{(x,y_1,y_0,g) \sim \mu_{\tilde{\mathcal{D}}}} \Big[D^2_{\text{Hellinger}} \big(\mathbb{P}_{R^*} (\cdot | x, y_1, y_0, g) \| \mathbb{P}_R (\cdot | x, y_1, y_0, g) \big) \Big] \\ + \frac{3}{N} \cdot \log \left(\frac{\mathcal{N}_{\varepsilon}(\mathcal{R}, \| \cdot \|_{\infty})}{\delta} \right), \tag{A.30}$

where $\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})$ denotes the ε -covering number (Zhou, 2002) of the reward model class \mathcal{R} . By the relationship between the Hellinger distance and TV distance, we have

1120
$$D^2_{\text{Hellinger}} \left(\mathbb{P}_{R^{\star}}(\cdot|x, y_1, y_0, g) \| \mathbb{P}_R(\cdot|x, y_1, y_0, g) \right) \ge D^2_{\text{TV}} \left(\mathbb{P}_{R^{\star}}(\cdot|x, y_1, y_0, g) \| \mathbb{P}_R(\cdot|x, y_1, y_0, g) \right),$$

By the definition of the goal-conditioned preference model \mathbb{P}_R in (A.1), we have

$$D_{\mathrm{TV}}(\mathbb{P}_{R^{\star}}(\cdot|x,y_{1},y_{0},g)||\mathbb{P}_{R}(\cdot|x,y_{1},y_{0},g))$$

$$=\frac{1}{2} \cdot \left|\sigma\left(R^{\star}(x,y_{1},g^{\star})-R^{\star}(x,y_{0},g^{\star})\right)-\sigma\left(R(x,y_{1},g^{\star})-R(x,y_{0},g^{\star})\right)\right|$$

$$+\frac{1}{2} \cdot \left|\sigma\left(R^{\star}(x,y_{0},g^{\star})-R^{\star}(x,y_{1},g^{\star})\right)-\sigma\left(R(x,y_{0},g^{\star})-R(x,y_{1},g^{\star})\right)\right|$$

$$=\left|\sigma\left(R^{\star}(x,y_{1},g^{\star})-R^{\star}(x,y_{0},g^{\star})\right)-\sigma\left(R(x,y_{1},g^{\star})-R(x,y_{0},g^{\star})\right)\right|, \quad (A.31)$$

where the second equality uses the fact that $\sigma(-z) = 1 - \sigma(z)$. Applying Lemma A.8 and the condition that $R(x, y, g) \in [B/2, B/2]$ for any $(x, y, R, g) \in \mathcal{X} \times \mathcal{A} \times \mathcal{R} \times \mathcal{G}$ in Assumption A.1, we have

$$\left|\sigma\left(R^{\star}(x,y_1,g^{\star})-R^{\star}(x,y_0,g^{\star})\right)-\sigma\left(R(x,y_1,g^{\star})-R(x,y_0,g^{\star})\right)\right|$$

1134
1135
$$\geq \kappa \cdot \left| \left(R^{\star}(x, y_1, g^{\star}) - R^{\star}(x, y_0, g^{\star}) \right) - \left(R(x, y_1, g^{\star}) - R(x, y_0, g^{\star}) \right) \right|, \quad (A.32)$$

where $\kappa = 1/(1 + \exp(B))^2$. Therefore, we bound the left-hand side of (A.33) as

 $\mathcal{L}(R^{\star}, \bar{\mathcal{D}}) - \mathcal{L}(R, \bar{\mathcal{D}})$

$$\leq -2\kappa^{2} \cdot \mathbb{E}_{(x,y_{1},y_{0},g)\sim\mu_{\tilde{\mathcal{D}}}}\left[\left|\left(R^{\star}(x,y_{1},g)-R^{\star}(x,y_{0},g)\right)-\left(R(x,y_{1},g)-R(x,y_{0},g)\right)\right|^{2}\right] + \frac{3}{N} \cdot \log\left(\frac{\mathcal{N}_{\varepsilon}(\mathcal{R},\|\cdot\|_{\infty})}{\delta}\right).$$
(A.33)

Meanwhile, the reward difference term in (A.29), which is evaluated on responses sampled from π^* and $\pi_{\rm sft}$, can be related to the reward difference evaluated on the data distribution $\mu_{\bar{D}}$ via Assump-tion A.3 as follows,

$$\mathbb{E}_{x \sim d_{0}, y_{1} \sim \pi^{\star}(\cdot|x, g^{\star}), y_{0} \sim \pi_{\mathrm{sft}}(\cdot|x, g^{\star})} \Big[\left(R^{\star}(x, y_{1}, g^{\star}) - R^{\star}(x, y_{0}, g^{\star}) \right) - \left(R(x, y_{1}, g^{\star}) - R(x, y_{0}, g^{\star}) \right) \Big] \\ \leq C_{\mu_{\overline{D}}}(\mathcal{R}; \pi^{\star}, \pi_{\mathrm{sft}}) \sqrt{\mathbb{E}_{(x, y_{1}, y_{0}, g) \sim \mu_{\overline{D}}} \left[\left| \left(R^{\star}(x, y_{1}, g) - R^{\star}(x, y_{0}, g) \right) - \left(R(x, y_{1}, g) - R(x, y_{0}, g) \right) \right|^{2} \right]}$$

$$(A.34)$$

Combining (A.33), (A.34), and (A.29) and defining

$$\Delta_R := \sqrt{\mathbb{E}_{(x,y_1,y_0,g)\sim\mu_{\bar{\mathcal{D}}}} \left[\left| \left(R^{\star}(x,y_1,g) - R^{\star}(x,y_0,g) \right) - \left(R(x,y_1,g) - R(x,y_0,g) \right) \right|^2 \right],$$
(A.35)

we obtain

$$\operatorname{Gap}^{\pi^{\star}}(\widehat{\pi}) \leq \max_{R \in \mathcal{R}} \left\{ C_{\mu_{\overline{D}}}(\mathcal{R}; \pi^{\star}, \pi_{\operatorname{sft}}) \cdot \Delta_{R} - 2\eta^{-1}\kappa^{2} \cdot \Delta_{R}^{2} \right\} + \frac{3}{\eta N} \cdot \log\left(\frac{\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})}{\delta}\right) \\ + \beta \cdot \mathbb{E}_{x \sim d_{0}} \left[\operatorname{KL}\left(\pi^{\star}(\cdot|x, g^{\star}) \|\pi_{\operatorname{ref}}(\cdot|x, g^{\star})\right) - \operatorname{KL}\left(\widehat{\pi}(\cdot|x, g^{\star}) \|\pi_{\operatorname{ref}}(\cdot|x, g^{\star})\right)\right] \\ \leq \frac{\left(C_{\mu_{\overline{D}}}(\mathcal{R}; \pi^{\star}, \pi_{\operatorname{sft}})\right)^{2} \eta}{8\kappa^{2}} + \frac{3}{\eta N} \cdot \log\left(\frac{\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})}{\delta}\right) \\ + \beta \cdot \mathbb{E}_{x \sim d_{0}} \left[\operatorname{KL}\left(\pi^{\star}(\cdot|x, g^{\star}) \|\pi_{\operatorname{ref}}(\cdot|x, g^{\star})\right)\right], \tag{A.36}$$

where the second inequality uses the fact that $az - bz^2 \leq a^2/(4b)$ for any $z \in \mathbb{R}$ and KL-divergence is non-negative. As a result, selecting $\varepsilon = (6 \cdot (1 + e^B)^{-1} \cdot N)^{-1}$ and

$$\eta = 2\sqrt{6} \cdot \sqrt{\frac{\log\left(\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})/\delta\right)}{N}}, \quad \beta = \frac{1}{\sqrt{N}}, \quad \kappa = \frac{1}{(1 + \exp(B))^2}, \tag{A.37}$$

we prove that with probability at least $1 - \delta$ that

$$\operatorname{Gap}^{\pi^{\star}}(\widehat{\pi}) \leq \sqrt{\frac{1}{N}} \cdot \left\{ \frac{\sqrt{6}}{4} \left(1 + \exp(B) \right)^{2} \left(\left(C_{\mu_{\widehat{D}}}(\mathcal{R}; \pi^{\star}, \pi_{\mathrm{sft}}) \right)^{2} + 1 \right) \iota \right. \\ \left. + \mathbb{E}_{x \sim d_{0}} \left[\operatorname{KL} \left(\pi^{\star}(\cdot | x, g^{\star}) \| \pi_{\mathrm{ref}}(\cdot | x, g^{\star}) \right) \right] \right\},$$
(A.38)

where we denote $\iota = \sqrt{\log (\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})/\delta)}$. Combining Theorem A.6, (A.20), and (A.38), we conclude the proof of Theorem A.4.

A.6 TECHNICAL LEMMAS

Lemma A.7 (Uniform concentration). Consider the negative log-likelihood loss in (A.2) and define the approximation error as $\varepsilon = (6 \cdot (1 + e^B) \cdot N)^{-1}$, where we assume that $R(x, y, q) \in [-B/2, B/2]$

1188 for any $(R, x, y, q) \in \mathcal{R} \times \mathcal{X} \times \mathcal{Y} \times \mathcal{G}$. Suppose that the reward model class \mathcal{R} has a finite ε -covering 1189 number $\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty}) < \infty$. Then for any $\delta < 1/e$ it holds with probability at least $1 - \delta$ that 1190 $\mathcal{L}(R^{\star}, \bar{\mathcal{D}}) - \mathcal{L}(R, \bar{\mathcal{D}})$ (A.39) 1191 $\leq -2 \cdot \mathbb{E}_{(x,y_1,y_0,g) \sim \mu_{\vec{D}}} \Big[D^2_{\text{Hellinger}} \big(\mathbb{P}_{R^\star}(\cdot | x, y_1, y_0, g) \big\| \mathbb{P}_R(\cdot | x, y_1, y_0, g) \big) \Big]$ 1192 1193 $+ \frac{3}{N} \cdot \log\left(\frac{\mathcal{N}_{\varepsilon}(\mathcal{R}, \|\cdot\|_{\infty})}{\delta}\right).$ 1194 (A.40)1195 1196 *Proof.* See the proof of Lemma D.1 in Liu et al. (2024), where we use the fact that (x, g) follows a 1197 fixed distribution. 1198 1199 **Lemma A.8** (Difference of Sigmoid functions). For any real numbers $z_1, z_2 \in [-B/2, B/2]$, it holds that 1201 $\kappa \cdot |z_1 - z_2| \le |\sigma(z_1) - \sigma(z_2)| \le |z_1 - z_2|,$ (A.41)1202 1203 where the constant $\kappa = 1/(1 + \exp(B))^2$. 1204 1205 *Proof.* See the proof of Lemma D.2 in Liu et al. (2024). 1206 1207 **EXPERIMENT DETAILS** 1208 В 1209 1210 **B**.1 SETUP 1211 We use the following prompt during training. Here, the reward values are the quality scores given by 1212 1213

we use the following prompt during training. Here, the reward values are the quarty scores given by
the judge models that exist in the preference dataset. The prompt is set as the system prompt whenever the LLM supports, such as Qwen2-7B-Instruct and Llama-3.1-8B-Instruct, and it is prefixed
before the original prompt when the LLM doesn't support system prompting, such as Mistral-7BInstruct-v0.3 and Gemma-2-9B-It.

```
Training prompt
You are an assistant that generates responses for the instruction
while implicitly achieving the following target score (on a scale of
1-10, where 1 is lowest and 10 is highest):
Overall score: {reward_value}.
```

At inference time, we use almost the same prompt, except that the goal score is the highest one, i.e., the overall score is 10.

Inference prompt

```
You are an assistant that generates responses for the instruction while implicitly achieving the following target score (on a scale of 1-10, where 1 is lowest and 10 is highest):
Overall score: 10.
```

1229 1230

1217

1218

1219

1220

1222

1223

1224

1225 1226

1227

1228

In our experiments using UltraFeedback, we directly leverage the LLM-as-Judge scores provided by GPT-4 in the dataset, which range from 1 to 10. For our method that is applied to on-policy data ranked by external reward models, including PairRM and ArmoRM, we apply linear transformations to normalize the resulting reward scores, ensuring they are scaled within the same 1 to 10 range.

For hyperparameters, we utilize a KL regularization coefficient of $\beta = 0.01$ in DPO, and we adopt the AdamW optimizer (Loshchilov, 2017). The batch size is set to 128, with a learning rate of 5e - 7and a warmup ratio of 0.1. Furthermore, we observe that for models such as Qwen2-7B-Instruct and Gemma-2-9B-It on UltraFeedback, as well as Llama-3-8B-Instruct on on-policy data, both DPO and our proposed method yield improved performance when employing the conservative DPO (cDPO) technique (Mitchell, 2023). Consequently, for these models, we set the label smoothing hyperparameter from the Alignment Handbook (Tunstall et al., 2023a) to 0.3, while keeping it at 0 for the remaining models.

1242 **B.2** FULL RESULTS 1243

1244	In Table 11, we present the full results on instruction-following benchmarks, which correspond to
1245	the performance illustrated in Figure 2 in the main text.

	AlpacaEval 2.0			MT-Bench			Arena-Hard-Auto	
	LC WR	WR	Avg. Len.	Avg.	1st	2nd	Score	Avg. Len.
Mistral-7B-Instruct-v0.3	19.65	15.40	1503	7.67	8.00	7.34	17.0	494
+DPO (UltraFeedback)	18.76	16.93	1643	7.66	7.92	7.40	17.6	504
+DPO (Reward-Augmented)	25.99	28.36	2270	7.69	8.02	7.36	18.3	883
Qwen2-7B-Instruct	20.93	18.22	1788	7.90	8.23	7.56	24.3	617
+DPO (UltraFeedback)	21.46	19.35	1797	8.33	8.72	7.93	21.9	553
+DPO (Reward-Augmented)	31.17	27.58	1789	8.47	8.93	7.97	30.1	644
Llama-3.1-8B-Instruct	24.79	27.38	2081	8.44	8.99	7.90	26.9	831
+DPO (UltraFeedback)	28.67	30.21	2053	8.47	9.01	7.93	33.0	1070
+DPO (Reward-Augmented)	31.20	35.93	2006	8.47	8.91	8.03	34.4	824
Gemma-2-9B-It	49.20	37.58	1572	8.54	8.81	8.28	42.8	541
+DPO (UltraFeedback)	50.70	35.02	1464	8.54	8.70	8.37	35.8	456
+DPO (Reward-Augmented)	59.27	54.56	1872	8.59	8.93	8.25	43.9	611
SPPO	55.60	49.61	1822	8.40	8.53	8.26	47.6	578
+DPO (UltraFeedback)	52.75	40.58	1544	8.41	8.78	8.04	40.4	457
+DPO (Reward-Augmented)	60.97	66.41	2543	8.73	9.06	8.41	49.0	761

Table 11: Results of the DPO models fine-tuned on UltraFeedback and on reward-augmented Ul-1262 traFeedback. We evaluate on the instruction-following benchmarks including AlpacaEval 2.0, MT-1263 Bench, and Arena-Hard-Auto. 1264

We also provide the full comparison results with reward-augmented methods in Table 12.

1268		Zephyr-SFT	DPO	DPA	SteerLM	NCA-P	NCA-R	INCA-P	INCA-R	Ours
1269	LC Win Rate	6.21	11.60	11.13	-	11.50	12.87	13.68	14.83	16.66
1270	Win Rate	3.94	8.58	10.58	8.21	8.43	9.56	11.00	11.34	13.37
1271	Avg. Len.	893	1240	1671	1585	1287	1364	1449	1338	1812

127 1272 1273

1285

1265 1266

1267

Table 12: Full comparison results with reward-augmented methods.

B.3 MORE ABLATIONS 1274

1275 **Controllable Generation with Prompt.** 1276

In Table 13, we ablate how generations 1277 differ when changing the goal rewards in 1278 the system prompt. We observe that the 1279 AlpacaEval 2.0 scores of the Qwen2-7B-1280 It+DPO (RA) model change accordingly 1281 as q varies. However, using the same 1282 q = 10 prompt during inference for the

	g = 10	g = 8	g = 6	UF(g = 10)
LC WR	31.17	28.66	25.56	24.44
WR	27.58	25.57	18.88	20.75

Table 13: Performance when conditioned on different goal rewards in the inference prompt.

Qwen2-7B-It+DPO (UF) model fails to give competitive results, indicating that our method is supe-1283 rior not only because of the additional system prompt. 1284

Benefits of Learning from High-Quality Rejected Responses. Using the UltraFeedback dataset, 1286 we construct two reward-augmented preference datasets by filtering out augmented data based on 1287 rejected responses with low and high reward values, respectively. Compared to our method, these 1288 datasets isolate the impact of excluding low- and high-reward rejected responses as goals. The 1289 evaluation results on AlpacaEval 2.0 are presented in Table 14. Learning from rejected high-reward 1290 samples demonstrates superior performance compared to the approach that excludes these samples. 1291

1292		0 0 TD 1			+DPO (RA	+DPO (RA
1293		Qwen2-7B-It	+DPO (UF)	+DPO (RA)	filter high)	filter low)
1294	LC Win Rate	20.93	21.46	31.17	29.36	31.81
1295	Win Rate	18.22	19.35	27.58	27.04	27.28

Table 14: Ablation on the benefits of learning from high-quality rejected responses.

Impact of the Reward Scale. For the UltraFeedback dataset that contains response rewards in the range of 1-10, we relabel them to be in the range of 1-5 and 1-100 with linear transformation. Our method followed by DPO is then applied on these different scaled datasets. The results are shown in Table 15. It can be observed that our method is robust to the reward scales. Since our main experiments use the default 1-10 scale as in UltraFeedback, it is likely that the performance can be further boosted, e.g., by adopting the 1-100 scale.

1303		Qwen2-7B-It	+DPO (UF)	+DPO (RA, 5)	+DPO (RA 10)	+DPO (RA 100)
1304	LC Win Rate	20.93	21.46	29.85	31.17	31.81
1305	Win Rate	18.22	19.35	26.12	27.58	27.96

Table 15: Ablation on the impact of the reward scale demonstrates the robustness of our method.