DATA-CENTRIC HUMAN PREFERENCE OPTIMIZATION WITH RATIONALES

Anonymous authors

003 004

010 011

012

013

014

015

016

017

018

019

021

025

026

Paper under double-blind review

Abstract

Reinforcement learning from human feedback plays a crucial role in aligning language models towards human preferences, traditionally represented through comparisons between pairs or sets of responses within a given context. While many studies have enhanced algorithmic techniques to optimize learning from such data, this work shifts focus to improving preference learning through a data-centric approach. Specifically, we propose enriching existing preference datasets with machine-generated rationales that explain the reasons behind choices. We develop a simple and principled framework to augment current preference learning methods with *rationale* information. Our comprehensive analysis highlights how rationales enhance learning efficiency. Extensive experiments reveal that rationale-enriched preference learning offers multiple advantages: it improves annotation efficiency, accelerates convergence to higher-performing models, and reduces verbosity bias and hallucination. Furthermore, this framework is versatile enough to integrate with various preference optimization algorithms. Overall, our findings highlight the potential of re-imagining data design for preference learning, demonstrating that even freely available machine-generated rationales can significantly boost performance across multiple dimensions.

028 1 INTRODUCTION

Preference tuning is an important step in the language model training process so as to productize and
 deploy them, whose goal is to align the model towards human preferences and prevent the model
 from unwanted behavior (Christiano et al., 2017; Stiennon et al., 2020; Bakker et al., 2022).

These preferences are typically introduced into the dataset through prompts and ranked responses. This dataset is then utilized by reinforcement learning from human feedback (RLHF) methods (Ouyang et al., 2022) to optimize reward or preference models, thereby aligning the language model. (Schulman et al., 2017) proposes a reinforcement learning-based algorithm to train the model by maximizing the reward given the reward function, where the first step involves learning a reward model to replicate human preferences. Alternatively, (Rafailov et al., 2024) proposes a direct preference optimization (DPO) algorithm, which avoids training a separate reward model and optimizes the policy through implicit Bradley-Terry reward modeling with a single objective.

However, these methods often face several challenges such as overfitting Azar et al. (2024), performance degradation Pal et al. (2024), reward exploitation Amodei et al. (2016), or the generation of excessively long inputs Park et al. (2024). In addition, collecting preference datasets can be costly Tan et al. (2024), but without sufficient samples, these methods would risk underfitting Jinnai & Honda (2024). Various studies aim to address these issues through improved algorithmic designs, either by regularizing the objective Pal et al. (2024); Amini et al. (2024); Park et al. (2024) or by introducing new formulations Ethayarajh et al. (2024); Hong et al. (2024); Yuan et al. (2023); Munos et al. (2023); Swamy et al. (2024); Wu et al. (2024).

In our study, we transition from an algorithmic to a data-centric approach, posing the key question:
 How can enhancing the preference dataset aid the model in boosting its performance and efficiency in preference learning? By rethinking preference learning through the lens of data, we aim to discover new insights and opportunities that can unlock the potential for more robust, data-efficient preference learning. Given the current setup of preference datasets, an important question arises: why would a certain response be preferred over another? For obvious cases, it is simple for humans to understand

054

056

058

060

061

062

063

064 065

066 067

068

069

071

090

092

095

096



Figure 1: Comparison between the current pair-wise preference dataset used for preference learning and the enriched dataset with added rationales.

the preference. However, when the responses are closely matched, understanding the preference without any explanation becomes challenging. Even superficial features such as length might not 073 serve as a straightforward metric to determine preference. In one instance, a longer response may be 074 favored for its comprehensiveness, while in another instance, a shorter response might be preferred for 075 its conciseness. Another consideration regarding preferences is that individuals might have varying 076 preferences for different reasons, and without explicitly outlining these reasons, one would be unable 077 to discern the underlying rationale. Given these challenges, the model will struggle to learn these 078 preferences without any explanations, causing data inefficiency, and worse, it could learn the wrong 079 cues, decreasing performance.

For the reasons outlined above, we propose a natural extension to the current dataset structure by enriching the preferences with rationales, aiding the model in better understanding during preference learning. Rationales explain why one response is preferred over another for a given prompt. This idea also draws inspiration from social studies (Mitchell et al., 1986; Chi et al., 1994; Crowley & Siegler, 1999; Williams & Lombrozo, 2010), showing that adding explanations to answers improves one's understanding of the problem compared to individuals who do not provide explanations.

Contributions. This paper provides a new data-centric perspective on preference learning. We list the summary of contributions:

- We introduce rationales into the human preference learning framework, where rationales explain the reasons behind the preference for a particular response. In practice, these rationales can be generated in a cost-effective manner by prompting an off-the-shelf language model, which may or may not have undergone preference learning.
 - We derive a straightforward formulation for the preference function to extend the rationales and show how to adapt our method to current preference learning algorithms such as DPO.
 We analytically examine the impact of the rationale on preference training through the lens of information theory. Our theoretical analysis demonstrates that highly informative rationales can improve preference prediction accuracy and reduce sample complexity.
- 097rationales can improve preference prediction accuracy and reduce sample complexity.098I We empirically show the impact of preference learning with rationales, highlighting improvements in both performance and annotation efficiency compared to baseline methods.
Specifically, the rationale-enriched DPO model can save up to $3 \times$ the annotated data required by the vanilla DPO model. With the same amount of data, it can improve the winrate against the supervised fine-tuned model by 8 9%. Further, the rationale-based DPO model shows reduced susceptibility to verbosity bias and truthfulness degradation compared to DPO. We demonstrate the flexibility and effectiveness of our approach by extending rationales to ORPO.
- 105I We showcase the efficacy of rationales generated by the off-the-shelf models with $\leq 8B$ 106parameters on the preference learning. We emphasize the importance of high-quality data107for improved preference learning.
 - I We release our code and datasets to facilitate further research in this direction.

In a broader context, our approach presents a new paradigm for data-centric research in language modeling: rather than focusing on pruning samples to distill the most informative pieces from a dataset Albalak et al. (2024), we explore how to enrich each sample's information content and examine its impact. The promising results presented in this paper demonstrate the effectiveness of enhancing individual samples' information content in preference learning and suggest that this approach may hold potential for improving learning in other domains.

114 115

116

2 RELATED WORK

117 **RLHF with Reward Modeling.** Tuning large language models to align their outputs towards human 118 preferences is crucial for controlling model behavior and maintaining desirable boundaries (Casper 119 et al., 2023). To achieve this, RLHF has been introduced; aligning models through preference 120 training (Christiano et al., 2017; Ziegler et al., 2019; Stiennon et al., 2020). Schulman et al. (2017) describe a method that typically involves two stages. The first stage learns a reward model using a 121 preference dataset often modeled under the Bradley-Terry model Bradley & Terry (1952). The second 122 stage fine-tunes the target policy model to maximize the rewards from the reward model, employing 123 algorithms such as proximal policy optimization (PPO) proposed by Schulman et al. (2017) and 124 adopted in Ouyang et al. (2022). A direct preference optimization (DPO) method that implicitly 125 models the reward function was introduced by Rafailov et al. (2024). However, Azar et al. (2024) 126 observe that RLHF and DPO are prone to overfitting due to the assumptions of the Bradley-Terry 127 model. Conversely, Pal et al. (2024) explore the possibility of DPO underfitting when dealing with 128 challenging responses that are difficult for the model to distinguish. Additionally, Park et al. (2024) 129 note that DPO can exploit response length to maximize reward, proposing a length-regularized DPO 130 (R-DPO) to address this issue. This should not be confused with our rationale-based DPO (RDPO) 131 method. Interestingly, we observe that if rationales mention conciseness as a feature, then the length of responses is significantly reduced compared to SFT and DPO responses. The learning dynamics 132 during preference tuning are analyzed in Im & Li (2024), emphasizing the importance of high-quality 133 preference datasets for effective learning. They find that the more distinguishable the response pairs, 134 the easier it is for the model to learn, leading to faster convergence. This has also been observed 135 in Pal et al. (2024). However, designing such datasets is challenging, and it also remains important 136 for models to learn from such nuanced responses, which appear in practice. We try to address the 137 difficulty of model learning from intricate preferences by providing rationales during preference 138 training. To improve efficiency over DPO, which requires an intermediate step to train the reference 139 model, odds ratio preference optimization (ORPO) was introduced by Hong et al. (2024) to eliminate 140 this step. Another method by Ethayarajh et al. (2024) adapts the Kahneman-Tversky human utility 141 model to handle preference datasets with a single response (either chosen or rejected), removing the 142 need for training the model on both responses. Conversely, Yuan et al. (2023) propose a preference method that considers multiple ranked responses for a prompt and optimizes over them. Our method 143 can complement these methods by adding rationales into training. In this paper, we demonstrate an 144 extension of our framework to ORPO. 145

- General Preference Modeling. Reward modeling, however, can incentivize undesirable behaviors, 146 such as "reward hacking" (Amodei et al., 2016), where agents maximize rewards without achieving the 147 desired objective. Overfitting is another challenge, as exemplified in Azar et al. (2024). While effective 148 for comparing two responses, the Bradley-Terry preference modeling relies on the assumption of 149 transitivity, which may not hold true in practice (Bertrand et al., 2023). To address this, Munos 150 et al. (2023) introduced general preference modeling, which directly learns general preferences 151 by formulating a two-player, constant-sum game between policies. The goal is to maximize the 152 probability of generating the preferred response against the opponent. The solution is the Nash equilibrium of this game, where payoffs are derived from the general preference function. Building 153 upon this work, Munos et al. (2023) proposed an algorithm for the regularized general preference 154 model, while Swamy et al. (2024) developed a solution for the unregularized formulation and 155 introduced self-play preference optimization (SPO) as an iterative algorithm to reach the optimal 156 solution. However, SPO suffers from data inefficiency due to its two-timescale update rules. To 157 address this, Rosset et al. (2024) introduced an efficient direct Nash optimization (DNO) method 158 that leverages the DPO formulation in practice. Additionally, Wu et al. (2024) proposed an efficient, 159 scalable, iterative self-play method that generates responses generally preferred over others. 160
- 161 While previous efforts have introduced algorithmic enhancements for preference tuning, they have been limited to the existing framework of preference datasets with prompts and ranked responses. In

162 contrast, our work is first to introduce rationales, a data-centric solution, into preference learning. 163 Learning with Rationales. The supervised learning framework typically involves training a model 164 to learn the ground truth label for a given prompt without providing explicit explanations for the 165 associations, which can lead to the model learning incorrect cues. To mitigate this issue, rationales 166 have been integrated into the framework, offering explanations for the given associations. These rationales initially were generated by humans (Zaidan et al., 2007; Ross et al., 2017; Hase & Bansal, 167 2021; Pruthi et al., 2022). However, due to the high cost of human labor and the development of more 168 capable large language models, rationales are now often automatically generated by these models, reducing the need for human involvement (Wei et al., 2022; Kojima et al., 2022). Given rationales, 170 they have been used as guiding aids by incorporating them directly into the prompt during the training 171 phase (Rajani et al., 2019; Zelikman et al., 2022; Huang et al., 2023) or at the inference stage (Wei 172 et al., 2022; Kojima et al., 2022; Wang et al., 2022). Besides using them as additional context within 173 the prompt, rationales can also serve as labels to train models to generate such explanations for their 174 predictions (Wiegreffe et al., 2021; Narang et al., 2020; Eisenstein et al.; Wang et al.; Ho et al., 2023; 175 Magister et al., 2022; Li et al., 2023a). In similar manner, rationales have been applied in knowledge 176 distillation, where they are generated by a more capable models to supervise weaker models (Hsieh 177 et al., 2023; Chen et al., 2024). In parallel with these advancements, we introduce rationales into the preference learning landscape, where rationales are used to explain the preference of one answer over 178 another. Our findings demonstrate the effectiveness of rationales in preference learning, even when 179 generated by the same model or a smaller-sized model. 180

Synthetic Preference Data Generation. Synthetic preference data generation plays a pivotal role in 181 preference learning by creating new annotated datasets that capture user choices or preferences Yang 182 et al.; Pace et al.; Meng et al. (2024). These methods focus on producing preference pairs that can serve 183 as training data for models, enabling the exploration of diverse scenarios and reducing reliance on 184 costly manual annotations. Furthermore, Wu et al. (2024); Wang et al. (2024) try to further synthesize 185 the preference examples by iteratively generating new response pairs. However, our approach diverges fundamentally from this objective. While synthetic data generation targets the creation of new datasets, 187 our work emphasizes enhancing existing datasets by incorporating rationales, thereby enriching 188 preference annotations with explanatory depth. This distinction highlights the complementary nature of these methods: data generation addresses the early stage of creating foundational datasets, whereas 189 rationale augmentation enhances the interpretability and utility of existing data. Attempting to 190 compare these approaches directly would obscure their unique contributions. Instead, their synergy 191 lies in how synthetic data generation can produce the raw preference pairs that are later refined through 192 rationale augmentation, advancing the overall preference learning pipeline. Exploring how these 193 two approaches can be combined to create synthetic datasets enriched with rationales is an exciting 194 direction for future work, holding promise to further enhance the capabilities and generalizability of 195 preference learning models. 196

197 198

199

200

201

202

203

204 205

3 Method

In this section, we introduce the incorporation of rationales into preference learning and show the derivation of adapting current methods. We present a demonstration of extending the direct preference optimization (DPO) algorithm to incorporate the rationales, while similar extensions can be applied to other variants of DPO. Further, we analyze theoretically the possible impact of rationales through the perspective of information theory.

3.1 PRELIMINARIES

Solutions. Let \mathcal{D} denote the pair-wise preference dataset of size N, $\mathcal{D} = \{x^{(i)}, y^{(i)}_w, y^{(i)}_l\}_{i=1}^N$, where $x^{(i)}$ is a context, $y^{(i)}_w$ is the preferred/chosen/winning response to the context $x^{(i)}$ over the unpreferred/rejected/losing response $y^{(i)}_l$. Let π_θ and π_{ref} denote the policy to be preference optimized and the reference policy respectively. In our setting, the policies are the language model to be preference trained and the base or supervised fine-tuned SFT model, respectively. To compute the joint probability of the autoregressive language model π generating the response y given the prompt x, we compute the product of probabilities after observing each token: $\pi(y|x) = \prod_{t=0}^{|y|} \pi(y_t|x, y_{0:t})$.

- 214
- **Reward Modeling with DPO.** In the RLHF process (Christiano et al., 2017; Ziegler et al., 2019; Stiennon et al., 2020; Bai et al., 2022; Ouyang et al., 2022; Rafailov et al., 2024), the goal is to align the

216 language model towards human preferences. The preferences ranking from the dataset \mathcal{D} is assumed 217 to be sampled from the latent reward function $r^*(x, y)$ and the preference function is assumed to be 218 generated by the Bradley-Terry model (Bradley & Terry, 1952): $p^*(y_w \succ y_l|x) = \sigma(r^*(x, y_w) - \sigma(r^*(x, y_w)))$ 219 $r^*(x, y_l)$, where σ is the sigmoid function. The reward function then can be estimated by minimizing the log-likelihood of the following objective $\mathcal{L}(r, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r(x, y_w) - r(x, y_l))]$. Then the next step is to tune the language model with the reward model as follows by maximizing 220 221 the rewards and not diverging from the fixed reference model: $\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} [r(x, y)] - \sum_{x \in \mathcal{D}, y \in \mathcalD, y \in$ 222 $\beta \mathbb{D}_{\mathrm{KL}} [\pi_{\theta}(y|x) \| \pi_{\mathrm{ref}}(y|x)],$ where β is a hyperparameter measuring the divergence between two 223 policies. Alternatively, with a reparametrization of the Bradley-Terry preference model (Rafailov 224 et al., 2024), the preference function can be expressed in terms of policy π^* : 225

$$p^*(y_w \succ y_l|x) = \sigma \left(\beta \log \frac{\pi^*(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi^*(y_l|x)}{\pi_{\text{ref}}(y_l|x)}\right).$$
(1)

Thus, to estimate the policy, Rafailov et al. (2024) proposes to directly minimize the log-likelihood of the following DPO loss: $\mathcal{L}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right].$

3.2 FORMULATION OF PREFERENCE LEARNING WITH RATIONALES

233 While preferences are modeled given the preferred and unpreferred responses, there are nuances 234 in the responses that are obscure for the model to comprehend and catch the differences between 235 them. Therefore, our goal is to help the model learn the preferences by providing guidance cues in 236 the preference tuning process, which we call rationales. Rationales explain why a given response is 237 preferred over the other response. For that reason, we extend the current preference learning with a data-centric technique to incorporate rationales and we term this the rationale-enriched preference 238 *function*, where the updated preference function is formulated as $p^*(y_w \succ y_l, r|x)$ and r is the 239 rationale from the updated dataset $\mathcal{D}' = \{x^{(i)}, y^{(i)}_w, y^{(i)}_l, r^{(i)}\}_{i=1}^N$. By the chain rule, we arrive at: 240 241

$$p^{*}(y_{w} \succ y_{l}, r|x) = p^{*}(y_{w} \succ y_{l}|x) \cdot p^{*}(r|x, y_{w} \succ y_{l}),$$
(2)

242 where the first term is the pair-wise preference term modeled in Section 3.1, and the second term is 243 the probability of the rationale r given the context x and the preference $y_w \succ y_l$. Given the policy π^* , 244 we can retrieve the probability of generating the rationale r given the context x and the preference 245 $y_w \succ y_l, \pi^*(r|x, y_w \succ y_l)$. Similarly when retrieving the probability of generating responses y_w and 246 y_l for the prompt x, which are given in the preference dataset \mathcal{D}' , we can also retrieve the probability of generating rationale r given x, y_w , and y_l , where $(x, y_w, y_l, r) \sim \mathcal{D}'$. In practice, we ask the policy 247 language model to explain why the response y_w is preferred over the response y_l for the prompt x and 248 retrieve the probability of generating the rationale r. Thus, $p^*(r|x, y_w \succ y_l) = \pi^*(r|x, y_w \succ y_l)$. 249

Adaptation to DPO Loss. After deriving the rationale-enriched preference learning function, we extend the DPO method to incorporate rationales. By substituting $p^*(y_w \succ y_l|x)$ from Equation 1 and $p^*(r|x, y_w \succ y_l)$ into Equation 2, we can express the rationale-enriched preference function in terms of an optimal policy π^* :

$$p^{*}(y_{w} \succ y_{l}, r|x) = \sigma \left(\beta \log \frac{\pi^{*}(y_{w}|x)}{\pi_{\text{ref}}(y_{w}|x)} - \beta \log \frac{\pi^{*}(y_{l}|x)}{\pi_{\text{ref}}(y_{l}|x)}\right) \pi^{*}(r|x, y_{w} \succ y_{l}).$$
(3)

We can optimize our policy π_{θ} through maximum likelihood using the following objective over the updated preference dataset $\mathcal{D}' = \{x^{(i)}, y^{(i)}_w, y^{(i)}_l, r^{(i)}\}_{i=1}^N$, which we term as rationale-DPO (RDPO):

$$\mathcal{L}_{\text{RDPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l, r) \sim \mathcal{D}'} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) + \gamma \log \pi_{\theta}(r | x, y_w \succ y_l) \right], \quad (4)$$

where γ is the added hyperparameter for weighting the impact of rationales on the loss.

4 EVALUATION

226 227

231

232

255 256

257

263 264 265

266

In this section, we evaluate the impact of rationales on preference learning. We conduct multiple
 experiments with two main goals in mind: (1) to understand how the added rationales affect the
 efficacy and efficiency of current preference learning algorithms, and (2) to determine the significance of rationale quality for effective learning.

270 4.1 EXPERIMENTAL SETUP 271

272 Datasets. For our analysis, we focus on two popular preference datasets: Orca DPO Pairs (Intel, 273 2024), which is a pairwise preference dataset version of Orca (Mukherjee et al., 2023), and a binarized UltraFeedback (Tunstall et al., 2023), which is a pair-wise version of UltraFeedack (Cui et al., 2023). 274 For each dataset, we take 12,000 samples as training set and 512 fixed samples as test set for winrate 275 evaluations. We generate rationales and add rationales to the current datasets. We refer readers 276 to Appendix B.6 for details on generating rationales. For evaluating hallucination, we adopt the TriviaQA dataset (Joshi et al., 2017) and use LM Evaluation Harness (Gao et al., 2023) code to 278 measure the exact match (EM) accuracy on the dataset. Given the test sets, we sample the responses 279 from models trained with preference learning methods and compare the performance between the 280 models by measuring the winrates between the corresponding responses. 281

Models. We investigate preference training on various large language models: Mistral-7B-v0.1, 282 Mistral-7B-Instruct-v0.2 (Jiang et al., 2023), Zephyr-7B-Beta(Tunstall et al., 2023), and Llama3-8B-283 Instruct (AI@Meta, 2024). We use GPT-40 (Achiam et al., 2023) as a judge to evaluate the responses 284 generated by the models and to retrieve the winrate scores. While in this section, we mainly study the 285 Mistral-7B-Instruct-v0.2 model (unless explicitly specified) with rationales generated by this model, we also provide full results on remaining models with ablation on hyperparameters in Appendix B.2. 286 Methods. In our experiments, we study the integration of rationales into preference learning 287 frameworks, such as DPO (Rafailov et al., 2024), which requires the SFT model for the reference 288 model, and ORPO (Hong et al., 2024), which does not. To ensure fair comparison between DPO 289 and RDPO, we fine-tune the base model with supervised fine-tuning (SFT) only using the chosen 290 responses from the preference dataset for a single epoch. We extend the code implementation from 291 human-aware loss functions (HALOs) repository (Ethayarajh et al., 2023) to adapt to our methodology 292 and borrow the hyperparameters for each of the above methods in our study. 293



Figure 2: Winrate comparison between the models trained with RDPO and DPO. Left: Winrate against the SFT model trained on the Orca dataset. Right: Winrate against the SFT model trained on the Ultrafeedback dataset. X-axis denotes the training data size used for preference training of DPO and RDPO models.

313 314 315

316

296

297

299

300

301

302

303

304 305

306

307

308

309

310

311

312

4.2 PERFORMANCE OF PREFERENCE LEARNING WITH RATIONALES

317 **Versus SFT.** We examine how adding rationales to the current preference learning algorithms 318 can impact performance. We compare the responses generated by the preference-aligned model 319 against the responses generated by the SFT model and measure the winrates scores using the GPT-40 320 evaluator. To study data annotation efficiency, we train the models on various training data sizes, 321 ranging from 1,000 to 12,000 data points, for both RDPO and DPO training. We observe on the left side of Figure 2 that both models, DPO and RDPO models, achieve a better winning rate over the 322 SFT model (more than 50%) on the Orca dataset with an increasing winning trend when training 323 data increases. Additionally, we observe as the DPO model converges to around 60% winning rate

340

341

342

343 344

345

346

347

348

349

350

351

352

353

354

355

356 357 358

359

360

361

362

363

364

366

324 at the 9,000 mark against the SFT model, the RDPO model achieves this rate at a smaller training 325 data size with 3,000 training data points. Furthermore, we observe that the RDPO model can reach 326 an even higher winning rate against the SFT than the DPO model, reaching above the 66% winning 327 rate. We see a similar observation with the models trained on the UltraFeedback dataset on the 328 right side of Figure 2. Furthermore, the drawing rate for the RDPO model is stable and low across different training data sizes, which shows that RDPO winning rate is higher not due to flipping 329 the draw points but the losing points. While RDPO can increase the computation time due to the 330 addition of rationales, the model trained on rationales can converge earlier with fewer data points 331 than DPO. This is especially important as the cost of collecting human preference data is high. Thus, 332 improving annotation efficiency can potentially save further training costs. Additionally, with enough 333 computation, RDPO can reach a better model than DPO does. 334

Moreover, we observe for the case of the UltraFeedback dataset, with more training data, the performance of the DPO model decreases. This can be attributed to the problem of DPO overfitting and exploiting length in longer responses Azar et al. (2024); Park et al. (2024). Indeed, the UltraFeedback dataset contains chosen responses that are longer on average (1,305 character length) than the rejected responses (1,124 character length). unlike the Orca dataset (784 and 978, respectively).



Figure 3: Winrate of RDPO model against the DPO model on respective datasets, Orca on the **left** and UltraFeedback on the **right**. The purple dashed line denotes the 0.5 mark.

Versus DPO. While we used SFT as a proxy to compare the performance between the DPO- and RDPO-trained models, here we directly compare the responses between these two models to measure the winrate for Orca and UltraFeedback datasets. We choose a DPO-trained model checkpoint for each dataset, where the winrate of the DPO model against the SFT model has converged, which is at 11,000 and 12,000, respectively. In Figure 3, we observe that the model trained with RDPO generates better responses on average than the model with DPO, even when trained with as little as 1,000 data points. With increasing training data, RDPO model improves the winrate to reach above 60% in both datasets. RDPO-trained model generates more of the preferred responses than the DPO-trained model does, even with $10 \times$ fewer training points.

367 **Response Comparison.** We compare the re-368 sponses generated by the DPO- and RDPO-369 trained models. As shown in Table 1, the av-370 erage output length by the DPO trained model 371 is much longer than the RDPO trained model in 372 the case of the Orca dataset, which is more than 373 5 times longer on average. Due to longer output, 374 there might be a chance for a higher occurrence 375 of hallucinations. Therefore, we want to study the correctness of the outputs from these mod-376 els. For this reason, we use the TriviaQA dataset 377 to measure the exact match (EM) accuracy and

	Avg Ou	tput Length	TriviaQA (Exact Match		
	DPO	RDPO	DPO	RDPO	
Orca UltraFeedback	2021 2066	364 1299	34.9 31.5	35.7 33.1	

Table 1: Comparison between DPO and RDPO. Left: The average output lengths of the generated responses on the prompts from the test sets of respective datasets. **Right:** The exact match (EM) performance on the TriviaQA dataset of the preference-trained models on respective datasets. compare between the models. We see in Table 1 that models trained with the DPO loss experience a
decrease in performance, compared to the models with RDPO loss. As a reason, we emphasize the
importance of measuring the hallucinations in the generations of both models in future studies. We
provide a comparison of the responses from the two models in Appendix B.9 and a time cost analysis
in Appendix B.4 to better guide model owners in method's tradeoff. We provide a comparison of the
responses from the two models in Appendix B.9 and a time cost analysis in Appendix B.4 to help
model owners better understand the trade-offs of our method.

386 Adaptation to ORPO. To demonstrate the flexibility of our rationale-enriched preference learning framework, we extend 387 the ORPO preference learning algorithm Hong et al. (2024), 388 which omits the SFT step, to include the rationales similar 389 to the RDPO loss, and we call it RORPO. As shown in Ta-390 ble 2, rationales can enhance the performance of ORPO and 391 achieve a better winrate over the vanilla ORPO trained model. 392 By successfully adapting rationales to both ORPO and DPO, 393 we emphasize the simplicity of the framework as well as the 394 effectiveness of rationales in preference learning. We further study the adaptation of these methods with rationales and evaluate the preference-trained models on the instruction-following 396

ORPO | 43 : 55 | RORPO

Table 2: Adapting rationales to the ORPO preference learning algorithm on Mistral-7B-v0.2-Instruct (Orca). Comparing the winrate of the ORPO- (**Left**) against the RORPO-trained model (**Right**).

benchmark, AlpacaEval 2.0 (Li et al., 2023b; Dubois et al., 2024), in Appendix B.4, and we observe consistent results when evaluating on this benchmark.

4.3 RATIONALE QUALITY ANALYSIS

385

399

400 401

402

403

404

421

422 423 In this section, we examine the importance of the quality of rationales for preference learning. We study different types of rationales and possible errors encountered in rationales, and how these affect the preference learning of the model.



Figure 4: Measuring the impact of types of rationales on the RDPO performance. Left: Comparing the winrates against the SFT model. **Right:** Comparing the winrates against the DPO model.

424 **Detailed Rationales vs General Rationales.** Explaining why one answer is preferred over the other 425 can be expressed in multiple ways through many perspectives. Here, we study the level of granularity 426 of the rationales, general (which explains the preference at a high level without going into details) 427 and detailed (which explains in details and pinpoints specifically to the prompt and the response). 428 We use language models to automatically generate the rationales for the Orca dataset according to our intent. For details on generation prompts, we refer to Appendix B.6. We provide samples of 429 these rationales in Appendix B.7. After training the models on respective rationales, we compare the 430 winrates between RDPO trained models and the DPO one. Figure 4 on the left shows that the model 431 trained on general rationales with the RDPO loss converges earlier to a high winning rate against the

432	Permuted Rationales V	S Original Rationales	Opposite Rationales VS Original Rationales				
433 434	<1	99	10	87			

Table 3: The analysis of quality of rationales on the RDPO performance. The winrate comparison
between the RDPO models trained on rationales with errors and original rationales. Left: Permuted,
irrelevant rationales. Right: Opposite, inaccurate rationales.

439 440

441

442

443

444

445

446

447

448

435

SFT model than the model trained on the detailed rationales. The reason could be that the general rationales share common features across the samples (e.g., clarity, conciseness, directness), which lets the model learn quickly and transfer these learning cues to other samples more easily, while detailed rationales might require more time to fully comprehend them. However, in both cases, the models trained on these rationales reach better winrates than the DPO against the SFT model. In Figure 4 (right), both RDPO models can have a better winrate > 57% against the DPO model with as few as 3,000 training samples, while the DPO model is trained on 11,000 samples. We provide results on models trained on additional epochs in Appendix B.3.

449 **Low-Quality Rationales.** RDPO has shown efficacy with rationales generated by the off-the-shelf 450 models, even when the models have not undergone preference alignment. However, we want to 451 further analyze the impact of rationale's quality on RDPO's performance. In particular, we examine the rationale quality in terms of its relevance and correctness. One case of a low-quality rationale can 452 be a completely irrelevant rationale to the given pair of responses. To simulate irrelevant rationales, 453 we permute the abovementioned detailed rationales over different samples so that no rationale is relevant to the context. Training the model on these rationales with RDPO and comparing one trained 455 on original rationales, we show in Table 3 that it achieves less than 1% winrate against the RDPO 456 model trained on correct and relevant rationales. To study the impact of correctness, we negate the 457 general rationales to have the opposite meaning and observe that the RDPO model trained on original 458 rationales gets almost a 90% winrate. As we note, the quality of rationales is important for improving 459 the preference learning performance. While we showed that rationales generated by the off-the-shelf 460 language models can already bring significant improvement to preference learning, we expect that 461 more deliberate control of the rationale quality can further improve preference learning. We leave the 462 in-depth exploration of strategies for generating quality-controlled rationales to future work.

Generated By	RDPO vs DPO Winrate				
Mistral-7B-Instruct-v0.2	76-23	50-46			
Llama3-8B-Instruct	73-27	52-45			
Phi3-Mini-4K	75-25	51-49			
	Mistral-7B-Instruct-v0.2	Llama3-8B-Instruct			
	Trained	On			

Table 4: Studying the impact of different source of rationale generation (**Y-axis**) for the Orca dataset on the model training with RDPO (**X-axis**). Winrate of the RDPO model against the DPO model.

472 473 474

Rationale Source. While collecting the human-annotated rationales could be high-quality, in 475 practice, it is costly with time and resources. Therefore, we resolve to language models to generate 476 rationales. In our experiments, we use the base models to create them. Here, we study the rationales 477 coming from other sources and how they impact the RDPO training. We generate rationales for 478 the Orca dataset on three different models: Mistral-7B-Instruct-v0.2, LLama3-8B-Instruct, and 479 Phi3-Mini-4K Abdin et al. (2024). Then, we use these rationales to train the first two models. Results 480 from Figure 4 show us consistent winrates against the DPO model with slightly better winrate from 481 the same source as the base model. This shows that the rationales can be transferred to other models 482 for preference training with rationales. Especially, leveraging models with small sizes $\{3, 7, 8\}$ billion parameters, we can generate rationales to improve preference learning. However, we observe 483 modest improvements on the Llama-3.1-8B-Instruct experiment, which can be attributed to two key 484 factors. First, the inherent capability of Llama-3.1-8B-Instruct surpasses Mistral-7B-Instruct-v0.2, 485 making substantial gains more challenging to achieve on this stronger baseline. Second, the general

preference datasets like Orca and Ultrafeedback, which include pre-existing responses, may not be fully optimized for Llama-3.1-8B-Instruct.

488 489 490

486

487

5 A SIMPLIFIED THEORETICAL MODEL FOR RATIONALES

To better understand rationale-enhanced preference learning, we employ machinaries from information theory to quantify benefits rationales provide in learning ground truth preferences under simplified assumptions. Formally, given query X, let the preference Z be a binary random variable, with Z = 1indicating that response Y_1 is preferred over response Y_2 , and Z = 0 indicating the opposite. Let R denote the rationale, $S = (X, Y_1, Y_2)$, and assume that the dataset $D = \{S_i, R_i, Z_i\}_{i=1}^n \sim \mu^n$ is sampled i.i.d. from a distribution μ .

As an intuitive first step, we quantify the benefits using the conditional mutual information between the true preferences and rationale-implied preferences given the input-response pairs, i.e., I(Z; g(R)|S), where g(R) capture the preference inferred from the rationale R. Intuitively, this mutual information quantity characterizes the added value of rationales in understanding true preferences, measuring how much additional information rationales provide beyond what can be inferred from input-response pairs S. Our analysis demonstrates a closed-form relationship between rationale informativeness and its alignment with true preferences (see Appendix A.1 for detailed derivations of all results).

504 Next, we quantify the benefits provided by rationales by directly analyzing the generalization error 505 for preference learning with and without rationales. Compared with analyzing the informativeness 506 of rationales outlined above, this approach offers a more direct measure of the impact of rationales 507 on learning outcomes. Note that the generalization error of preference learning still differs from the 508 winrate metrics used in our experiemnts for evaluating the generation of preference-aligned models. 509 However, they are intuitively connected under the DPO loss: the minimized test loss also indicates 510 learning the optimal data generation policy that achieves the highest reward. We defer the detailed 511 statement of the theorem and proof to Appendix A.2. In particular, our theoretical derivation shows that training with rationales can lead to improved lower generalization error when the rationale does 512 not contain irrelevant information other than those predictive of the preference Z, and the learning 513 process only captures the rationale information that is useful to predict Z. Despite being derived from 514 a simplified model, our theoretical insights align well with experimental observations. For instance, 515 our evaluation in Section 4.3 demonstrates that detailed rationales, which may contain extraneous 516 information, achieve lower sample efficiency compared to more general rationales that focus on 517 preference-predictive content; furthermore, when we manually inject noise into rationales to misalign 518 them with true preferences, we observe hampered preference learning, as indicated by lower winrates.

519 520 521

6 CONCLUSION & LIMITATIONS

522 In our work, we propose a paradigm shift in preference learning, emphasizing a data-centric per-523 spective. Language models are trained by presenting them with pairs of answers, one preferred and 524 one dispreferred, with the objective of teaching them human preferences. However, the selection 525 of preferred answers can often be ambiguous without explanation. To address this challenge and 526 enhance preference optimization efficiency, we introduce rationales that provide explicit reasoning 527 for choosing one answer over another. We propose a straightforward adaptation of existing losses by 528 incorporating these rationales into the training pipeline. Through extensive empirical experiments, 529 we demonstrate that rationales significantly enhance training data annotation efficiency and lead 530 to improved performance compared to baseline methods. Moreover, our approach is grounded in 531 information theory, offering insights into how rationales enhance preference training efficiency.

To date, we have integrated rationales into our training process and successfully trained models with up to 8 billion parameters using a dataset of 12,000 samples. We encourage further investigation into the impact of rationales on preference learning, particularly exploring larger models and datasets. To facilitate research in this area, we make our code and datasets publicly available.

With the development of unpaired preference learning algorithms, such as KTO Ethayarajh et al. (2024), it is important to extend the use of rationales to handle unpaired responses in future work, e.g., such as provided in the UltraFeedback dataset Cui et al. (2023) contains rationales for single responses without pairwise comparison.

540 REFERENCES 541

542 543 544	Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. <i>arXiv preprint arXiv:2404.14219</i> , 2024.
545 546 547	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> , 2023.
548 549 550	AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/ blob/main/MODEL_CARD.md.
551 552 553	Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, et al. A survey on data selection for language models. <i>arXiv preprint arXiv:2402.16827</i> , 2024.
554 555 556	Afra Amini, Tim Vieira, and Ryan Cotterell. Direct preference optimization with an offset. <i>arXiv</i> preprint arXiv:2402.10571, 2024.
557 558	Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. <i>arXiv preprint arXiv:1606.06565</i> , 2016.
559 560 561 562 563	Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 4447–4455. PMLR, 2024.
564 565 566	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. <i>arXiv preprint arXiv:2204.05862</i> , 2022.
567 568 569 570	Michiel Bakker, Martin Chadwick, Hannah Sheahan, Michael Tessler, Lucy Campbell-Gillingham, Jan Balaguer, Nat McAleese, Amelia Glaese, John Aslanides, Matt Botvinick, et al. Fine-tuning language models to find agreement among humans with diverse preferences. <i>Advances in Neural Information Processing Systems</i> , 35:38176–38189, 2022.
572 573 574	Quentin Bertrand, Wojciech Marian Czarnecki, and Gauthier Gidel. On the limitations of the elo, real-world games are transitive, not additive. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 2905–2921. PMLR, 2023.
575 576	Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. <i>Biometrika</i> , 39(3/4):324–345, 1952.
577 578 579 580 581	Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. <i>arXiv preprint arXiv:2307.15217</i> , 2023.
582 583	Xin Chen, Hanxian Huang, Yanjun Gao, Yi Wang, Jishen Zhao, and Ke Ding. Learning to maximize mutual information for chain-of-thought distillation. <i>arXiv preprint arXiv:2403.03348</i> , 2024.
584 585 586	Michelene TH Chi, Nicholas De Leeuw, Mei-Hung Chiu, and Christian LaVancher. Eliciting self-explanations improves understanding. <i>Cognitive science</i> , 18(3):439–477, 1994.
587 588 589	Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. <i>Advances in neural information processing systems</i> , 30, 2017.
590 591	Kevin Crowley and Robert S Siegler. Explanation and generalization in young children's strategy learning. <i>Child development</i> , 70(2):304–316, 1999.
592	Gangu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhivuan Liu.

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, 2023.

622

638

- 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.
- Jacob Eisenstein, Daniel Andor, Bernd Bohnet, Michael Collins, and David Mimno. Honest students
 from untrusted teachers: Learning an interpretable question-answering pipeline from a pretrained
 language model. In *Workshop on Trustworthy and Socially Responsible Machine Learning, NeurIPS 2022.*
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Model alignment as prospect theoretic optimization. In *Forty-first International Conference on Machine Learning*.
- 605 Kawin Ethayarajh, Winnie Xu, Dan Jurafsky, and Douwe Kiela. Human-2023. 606 aware loss functions (halos). Technical report, Contextual AI, https://github.com/ContextualAI/HALOs/blob/main/assets/report.pdf. 607
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 12 2023. URL https://zenodo.org/records/10256836.
- Peter Hase and Mohit Bansal. When can models learn from explanations? a formal framework for understanding the roles of explanation data. *arXiv preprint arXiv:2102.02201*, 2021.
- ⁶¹⁹ Namgyu Ho, Laura Schmid, and Se-Young Yun. Large language models are reasoning teachers. In
 Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 14852–14882, 2023.
- Jiwoo Hong, Noah Lee, and James Thorne. Reference-free monolithic preference optimization with
 odds ratio. *arXiv preprint arXiv:2403.07691*, 2024.
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 8003–8017, 2023.
- Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large
 language models can self-improve. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1051–1068, 2023.
- Shawn Im and Yixuan Li. Understanding the learning dynamics of alignment with human feedback.
 arXiv preprint arXiv:2403.18742, 2024.
- Intel. Intel. https://huggingface.co/Intel/neural-chat-7b-v3-1, 2024. Accessed: 2024-05-18.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- Yuu Jinnai and Ukyo Honda. Annotation-efficient preference optimization for language model alignment. *arXiv preprint arXiv:2405.13541*, 2024.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1601–1611, 2017.

648 649 650	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35: 22199–22213, 2022.
651	
652	Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. Symbolic
653	61st Annual Magting of the Association for Computational Linguistics (Volume 1: Long Papers)
654	nn 2665–2679 2023a
655	pp. 2005–2019, 2023a.
656	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy
657	Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following
658	models. https://github.com/tatsu-lab/alpaca_eval,2023b.
660	Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adamek, Eric Malmi, and Aliaksei Severyn,
661	Teaching small language models to reason. arXiv preprint arXiv:2212.08410, 2022.
662	Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-
663	free reward. arXiv preprint arXiv:2405.14734, 2024.
664	Tom M Mitchell Dishend M Kellen and Grander T Kalan Caladi'. Production has descent in the
665 666	A unifying view. <i>Machine learning</i> , 1:47–80, 1986.
667	Subhabrata Mukheriee Arindam Mitra Ganesh Jawahar Sahai Agarwal Hamid Palangi and Ahmed
668	Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. <i>arXiv preprint</i>
669	arXiv:2306.02707, 2023.
670	
671	Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland,
672	Znaonan Daniel Guo, Yunnao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, et al. Nash learning from human feedback arXiv praprint arXiv:2312.00886, 2023
673	carning from numan feedback. <i>urxiv preprint urxiv.2512.00000, 2025</i> .
675	Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma Malkan.
676	Wt5?! training text-to-text models to explain their predictions. arXiv preprint arXiv:2004.14546,
677	2020.
678	Long Quyang Jeffrey Wu Xu Jiang Diogo Almeida Carroll Wainwright Pamela Mishkin Chong
679	Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
680	instructions with human feedback. Advances in neural information processing systems, 35:27730-
681	27744, 2022.
682	Alizia Deep Jonethan Mallingon Eric Malmi Schootion Krouse and Aliekasi Severyn West of n
683	Synthetic preference generation for improved reward modeling. In ICLR 2024 Workshop on
684	Navigating and Addressing Data Problems for Foundation Models.
685	
686	Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White.
687	Smaug: Fixing failure modes of preference optimisation with dpo-positive. arXiv preprint
688	arxiv:2402.13228, 2024.
689	Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality in
690	direct preference optimization. arXiv preprint arXiv:2403.19159, 2024.
691	
692	Danish Pruthi, Rachit Bansal, Bhuwan Dhingra, Livio Baldini Soares, Michael Collins, Zachary C
693	tions from the teacher aid students? Transactions of the Association for Computational Linguistics
094 605	10:359–375, 2022.
696	
697	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
698	Finn. Direct preference optimization: Your language model is secretly a reward model. Advances
699	in iveural Information Processing Systems, 36, 2024.
700	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. Explain vourself!
701	leveraging language models for commonsense reasoning. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pp. 4932–4942, 2019.

702 703 704	Andrew Slavin Ross, Michael C Hughes, and Finale Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. <i>arXiv preprint arXiv:1703.03717</i> , 2017.
705 706 707	Corby Rosset, Ching-An Cheng, Arindam Mitra, Michael Santacroce, Ahmed Awadallah, and Tengyang Xie. Direct nash optimization: Teaching language models to self-improve with general preferences. <i>arXiv preprint arXiv:2404.03715</i> , 2024.
708 709 710	Daniel Russo and James Zou. Controlling bias in adaptive data analysis using information theory. In <i>Artificial Intelligence and Statistics</i> , pp. 1232–1240. PMLR, 2016.
711 712	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
713 714 715 716	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021, 2020.
717 718 719	Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A minimaxi- malist approach to reinforcement learning from human feedback. <i>arXiv preprint arXiv:2401.04056</i> , 2024.
720 721 722 723	Zhen Tan, Alimohammad Beigi, Song Wang, Ruocheng Guo, Amrita Bhattacharjee, Bohan Jiang, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. Large language models for data annotation: A survey. <i>arXiv preprint arXiv:2402.13446</i> , 2024.
724 725 726 727	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023.
728 729 730 731	PeiFeng Wang, Aaron Chan, Filip Ilievski, Muhao Chen, and Xiang Ren. Pinto: Faithful language reasoning using prompt-generated rationales. In <i>The Eleventh International Conference on Learning Representations</i> .
732 733 734	Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu, Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught evaluators. <i>CoRR</i> , 2024.
735 736 737 738	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh- ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> , 2022.
739 740 741	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837, 2022.
742 743 744 745	Sarah Wiegreffe, Ana Marasović, and Noah A Smith. Measuring association between labels and free-text rationales. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pp. 10266–10284, 2021.
746 747	Joseph J Williams and Tania Lombrozo. The role of explanation in discovery and generalization: Evidence from category learning. <i>Cognitive science</i> , 34(5):776–806, 2010.
749 750	Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play preference optimization for language model alignment. <i>arXiv preprint arXiv:2405.00675</i> , 2024.
751 752 753	Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. <i>Advances in neural information processing systems</i> , 30, 2017.
754 755	Kevin Yang, Dan Klein, Asli Celikyilmaz, Nanyun Peng, and Yuandong Tian. Rlcd: Reinforcement learning from contrastive distillation for lm alignment. In <i>The Twelfth International Conference on Learning Representations</i> .

756 757 758	Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback without tears. <i>arXiv preprint</i> <i>arXiv:2304.05302</i> , 2023.
759	Omar Zaidan Jason Fisner and Christine Piatko, Using "annotator rationales" to improve machine
760	learning for text categorization. In Human language technologies 2007: The conference of the
761	North American chapter of the association for computational linguistics: proceedings of the main
762	<i>conference</i> , pp. 260–267, 2007.
763	
764	Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with
765	reasoning. Advances in Neural Information Processing Systems, 35:154/6–15488, 2022.
700	Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
769	Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. arXiv
760	preprint arXiv:1909.08593, 2019.
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
790	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

810 A THEORETICAL DERIVATIONS

We begin with defining several standard quantities to be used throughout this section.

814 **Definition 1** Let X, Y and Z be arbitrary random variables, and let D_{KL} denote the KL divergence. 815 We denote P_X as the marginal probability distribution of X, and $P_{Y|X}$ as the conditional distribution. 816 The entropy of X is given by:

$$H(X) = -\sum_{x} P(X = x) \log P(X = x).$$

If X is a binary variable with p = P(X = 1) = 1 - P(X = 0), then we use H(p) for H(X).

The joint entropy of two random variables, H(X,Y), is the entropy of their joint distribution.

The conditional entropy of X given Y, H(X|Y), is:

$$H(X|Y) = H(X,Y) - H(Y)$$

The mutual information between X and Y is:

 $I(X;Y) = D_{KL}(P_{X,Y} || P_X P_Y).$

The disintegrated mutual information between X and Y given Z is:

 $I^{Z}(X;Y) = D_{KL}(P_{X,Y|Z} || P_{X|Z} P_{Y|Z}).$

The corresponding conditional mutual information is given by:

$$I(X;Y|Z) = \mathbb{E}_Z[I^Z(X;Y)].$$

If all entropies involved are finite, it can be shown that I(X;Y) = H(Y) - H(Y|X).

A.1 MUTUAL INFORMATION ANALYSIS

Formally, given query X, let the preference Z is a binary random variable, with Z = 1 indicating that response Y1 is preferred over response Y2, and Z = 0 indicating the opposite. Assume that the rationale-implied preference R is a binary random variable, with R = 1 indicating a rationale that supports Y1 being preferred, and R = 0 otherwise. For example, if the rationale mentions that Y1 is more concise and informative than Y2, then R = 1. However, there can be cases where $R \neq Z$, as the rationale may not always align perfectly with the actual preference.

847 For the analysis, we consider the following model: 1) The rationale R depends on the query-response 848 (QR) pair S = (X, Y1, Y2) and the preference Z, and is characterized by parameters β and α , where: $\beta = P(R = 1 | Z = 1, S) = P(R = 0 | Z = 0, S)$ represents the precision rate of consistency, and 849 $\alpha = P(R = 1 | Z = 0, S) = P(R = 0 | Z = 1, S)$ represents the recall error due to inconsistency. 2) 850 The preference Z is modeled as $P(Z = 1|S) = f(S) + \epsilon$, where f(S) captures the preference based 851 on the observable query-response pair S, and the additive noise term ϵ is a simple way to account for 852 unobserved factors influencing the complex human preference that are not captured in S. The term ϵ 853 is referred to as "bias" in the **main text** that accounts for the difference between the true preference 854 and prediction based on the query and responses alone. 855

Theorem 1 Under the given assumptions, the mutual information I(Z; R|S) is given by:

$$H(p+\epsilon) - (\beta(p+\epsilon) + \alpha(1-p-\epsilon)) \cdot H\left(\frac{\beta(p+\epsilon)}{\beta(p+\epsilon) + \alpha(1-p-\epsilon)}\right) \\ - (1 - (\beta(p+\epsilon) + \alpha(1-p-\epsilon))) \cdot H\left(\frac{\alpha(p+\epsilon)}{\alpha(p+\epsilon) + \beta(1-p-\epsilon)}\right)$$

861 862

856

817 818 819

821 822

823

824 825

826

827 828

829 830

831

832 833

838 839

840

where p = f(S). The mutual information I(Z; R|S) satisfies the following properties in three distinct regimes:

- 1. Uninformative rationale regime: If $\beta = \alpha = 0.5$, then I(Z; R|S) = 0.
- 2. Maximally informative rationale regime: If $\beta = 1$ and $\alpha = 0$, then $I(Z; R|S) = H(p + \epsilon)$.

867 3. Moderately informative rationale regime: If $\beta = 0.5 + \gamma$ and $\alpha = 0.5 - \gamma$, where $0 < \gamma < 0.5$, 868 then I(Z; R|S) increases with γ , ranging from 0 when $\gamma = 0$ (uninformative rationale) to $H(p + \epsilon)$ 869 when $\gamma = 0.5$ (maximally informative rationale).

The theorem highlights that the potential benefits of including rationales in the training process for preference learning tasks may vary in different regimes.

873 **Regime 1: Highly informative rationale** ($\beta = 1$ and $\alpha = 0$): In this case, the mutual information is solely determined by the entropy of the preference prediction from the query and responses, 874 $H(f(S) + \epsilon)$. Let us interpret f(S) as the query-response-dependent (QR-dependent) confidence 875 generator that only depends on S, and ϵ captures the idiosyncrasies such as unknown confounding 876 factors that influence the probability of preference. If the true probability $P(Z = 1|S) = f(S) + \epsilon$ 877 is less than 0.5, i.e., Z is most likely to be 0, a positive $\epsilon > 0$ means that the true probability 878 P(Z = 1|S) is less extreme than the confidence score f(S) as it gets closer to 0.5 with increasing ϵ . 879 On the contrary, a negative $\epsilon < 0$ means that the true probability P(Z = 1|S) is more extreme than 880 the QR-dependent confidence f(S) as it gets closer to 0 with increasing $|\epsilon|$. 881

From the perspective of the QR-dependent confidence generator f(S) that tries to explain the preference based on QR pairs, a positive $\epsilon > 0$ would make it look more confident than it should be, i.e., overconfident, while a negative $\epsilon < 0$ would make it less confident (or more conservative) based on the QR sequence than it should be.

If it is likely to be overconfident based on the QR pairs, i.e., $\epsilon > 0$, then the more positive ϵ is, the more risk there is of being overconfident in the QR-dependent predictor f(S). In this case, there is a lot of mutual information in I(Z; R|S), so having rationales can "soften" the potential overconfidence by bringing additional information other than the QR pair, which would otherwise occur in QR pairs alone, as in traditional reward modeling. Similar analysis holds when $P(Z = 1|S) = f(S) + \epsilon$ is greater than 0.5, i.e., Z is most likely to be 1.

Key message: Rationales are most useful when the reward modeling based on QR alone tends to have bias (i.e., overconfident).

Regime 2: Uninformative rationale ($\beta = \alpha = 0.5$): In this regime, the rationale provides no *additional* information about the preference, and the mutual information I(Z; R|S) is zero.

Regime 3: Moderately informative rationale (high precision $\beta = 0.5 + \gamma$ **and low recall error** $\alpha = 0.5 - \gamma$, where $0 < \gamma < 0.5$): In this regime, it can be shown based on dereivative analysis that as γ increases (more informative rationale), the terms involving γ in the numerators and denominators of the conditional entropies become more prominent. The mutual information will increase with γ , as the rationale becomes more informative about the preference.

901 902 903

A.2 THEOREM 2: GENERALIZATION BOUND

904 Next, we analyze the sample complexity of training language models with and without rationales to predict preferences. We consider two regimes: 1) Training with rationale: Let $\theta_{ra} = A_{ra}(D) \sim$ 905 $P_{\theta_{\pi}|D}$ denote the parameters of the language model trained to predict Z given S and R. 2) **Training** 906 without rationale: Let $\theta_{un} = A_{un}(D_{\backslash R}) \sim P_{\theta_{un}|D_{\backslash R}}$ denote the output parameters trained to 907 predict Z given only S, where $D_{\setminus R}$ is a dataset D with rationales removed. Given a loss function 908 ℓ that measures the prediction of preference Z, the (mean) generalization error is gen (μ, \mathcal{A}) = 909 $\mathbb{E}_{D,\theta \sim \mathcal{A}(D)} |\mathbb{E}_{\mu}[\ell(\theta)] - \mathbb{E}_{D}[\ell(\theta)]|$, where $\mathbb{E}_{\mu}[\ell(\theta)]$ is the expected loss on the true distribution (true 910 risk) and $\mathbb{E}_D[\ell(\theta)]$ is the empirical risk. 911

912 We introduce the following conditions on the relationship between S, R, and Z, and the learning 913 process: 1) $H(R|Z) \le \eta_1$ and $H(Z|R) \le \eta_2$, i.e., the rationale R is informative about Z (small η_2) 914 without excessive irrelevance (small η_1). 2) $I(\theta_{ra}; S|Z, R) \le \delta$ for some small positive constant δ , 915 i.e., the learned model θ_{ra} does not capture much additional information from S beyond what Z and 916 R already provide. Condition 1 is supported by an effective procedure to generate useful rationale R. 917 To justify Condition 2, if the learning algorithm is designed to focus on capturing the information in R, which is highly informative about Z (per Condition 1 for small η_2), we can show that the model 918 θ_{ra} can accurately predict Z without *needing* to capture much additional information from S beyond 919 what is already present in Z and R. We provide rigorous but partial justification in Appendix A.3.

Theorem 2 (Generalization bounds) Suppose the loss function ℓ is σ -subgaussian under the true data distribution. Under conditions 1 and 2, we have:

$$gen(\mu, \mathcal{A}_{ra}) \le \sqrt{\frac{2\sigma^2}{n} \cdot (I(\theta_{ra}; Z) + \delta + \eta_1)},$$
(5)

921 922

923 924

 $\operatorname{gen}(\mu, \mathcal{A}_{un}) \leq \sqrt{\frac{2\sigma^2}{n} \cdot (I(\theta_{un}; Z) + I(\theta_{un}; S|Z))}.$ (6)

The proof relies on the mutual information-based generalization bounds (Russo & Zou, 2016; Xu & 929 Raginsky, 2017) and the decomposition of the mutual information terms for both training regimes 930 using the chain rule (see Appendix A.2). The terms $I(\theta_{ra}; Z)$ and $I(\theta_{un}; Z)$ can be expected to be 931 similar as long as both regimes achieve good prediction of Z. Under the conditions of the theorem, 932 we can observe that the sample complexity reduction depends on the gap between $I(\theta_{un}; S|Z)$ and 933 $\delta + \eta_1$; training with rationales can lead to improved sample efficiency when the rationale does not 934 contain irrelevant information other than those predictive of the preference Z, i.e., η_1 is small, and the 935 learning process only captures the rationale information that is useful to predict Z. The theoretical 936 insights are supported by our experimental results. For instance, our evaluation in Section 4.3 937 demonstrates that a detailed rationale achieves lower sample efficiency compared to more general rationales, which contain less irrelevant information beyond what is predictive of the preference; 938 furthermore, we showed that irrelevant rationales, i.e., a large value of η_1 , indeed hamper learning. 939

940 For the proof, recall that given a loss function ℓ , the (mean) generalization error is gen (μ, \mathcal{A}) = 941 $\mathbb{E}_{D,\theta\sim\mathcal{A}(D)}|\mathbb{E}_{\mu}[\ell(\theta)] - \mathbb{E}_{D}[\ell(\theta)]|$, where $\mathbb{E}_{\mu}[\ell(\theta)]$ is the expected loss on the true distribution (true 942 risk) and $\mathbb{E}_D[\ell(\theta)]$ is the empirical risk. For fair comparison between gen (μ, \mathcal{A}_{un}) and gen $(\mu, \mathcal{A})_{ra}$, 943 some technical nuances arise. The key difference from the typical setup in Xu & Raginsky (2017) 944 is that the true data distribution μ includes the distribution for the rationale R, but the training 945 regime gen (μ, A_{un}) does not explicitly use this information. However, it may seem unclear initially whether we should include that in the generalization bound, since R is indeed generated based on Z, 946 corresponding to the true Markov chain: $S \rightarrow Z \rightarrow R$. 947

⁹⁴⁸ To clarify, the Markov chain for the training without rationale is: $S \to Z \to R$, with additional ⁹⁴⁹ arrows $Z \to \theta_{un}$ and $S \to \theta_{un}$, but no arrow from R to θ_{un} .

Intuitively, we should account for this difference by arguing that, conditioned on the preference Z, the learned model θ_{un} is conditionally independent of R. However, due to this difference, it seems prudent to reason from first principles.

Let's start by choosing the distributions P and Q for the Donsker-Varadhan variational representation of the KL divergence. We set $P = P_{S,R,Z,\theta_{un}}$ and $Q = \mu^n \otimes P_{\theta_{un}}$, where μ is the distribution for (S, R, Z). Then, for any measurable function f, we have:

$$D(P||Q) \ge \mathbb{E}_P[f(S, R, Z, \theta)] - \log \mathbb{E}_{(\bar{D}, \bar{R}, \bar{Z}, \bar{\theta}) \sim Q}[e^{f(S, R, Z, \theta)}].$$

$$\tag{7}$$

Now, choose $f(S, R, Z, \theta) = \lambda(\ell_D(\theta) - \ell_\mu(\theta))$ for some $\lambda \in \mathbb{R}$, where $\ell_D(\theta)$ is the empirical loss on the dataset D and $\ell_\mu(\theta)$ is the expected loss under the true distribution μ . Substituting this into equation 7, we get:

$$D(P||Q) \ge \lambda(\mathbb{E}[\ell_D(\theta)] - \mathbb{E}[\ell_\mu(\theta)]) - \log \mathbb{E}_{(\bar{D},\bar{R},\bar{Z},\bar{\theta})\sim Q}[e^{\lambda(\ell_{\bar{D}}(\theta) - \ell_\mu(\theta))}]$$

$$\ge \lambda(\mathbb{E}[\ell_D(\theta)] - \mathbb{E}[\ell_\mu(\theta)]) - \frac{\lambda^2 \sigma^2}{2n},$$
(8)

where the second inequality follows from the fact that $\ell_{\bar{D}}(\bar{\theta})$ is σ/\sqrt{n} -subgaussian under Q, due to the subgaussian assumption on the loss function.

As equation 8 holds for any $\lambda \in \mathbb{R}$, it must also hold for the λ that minimizes the right-hand side. This minimum occurs at $\lambda^* = n(\mathbb{E}[\ell_D(\theta)] - \mathbb{E}[\ell_\mu(\theta)])/\sigma^2$, yielding:

957 958

962 963 964

965

 $D(P||Q) \ge \frac{n}{2\sigma^2} (\mathbb{E}[\ell_D(\theta)] - \mathbb{E}[\ell_\mu(\theta)])^2.$

Taking the square root of both sides, we have:

$$\operatorname{gen}(\mu, \mathcal{A}_{\operatorname{un}}) \leq \sqrt{\frac{2\sigma^2}{n}D(P\|Q)}$$

976 The key observation here is that $P_{S,R,Z,\theta_{un}} = P_{\theta_{un}|S,R,Z}P_{S,R,Z} = P_{\theta_{un}|S,Z}P_{S,R,Z}$, since θ_{un} is 977 conditionally independent of R given S and Z in this training regime. Therefore,

$$D_{\mathrm{KL}}(P \| Q) = D_{\mathrm{KL}}(P_{S,R,Z,\theta_{\mathrm{un}}} \| \mu^n \otimes P_{\theta_{\mathrm{un}}})$$

= $D_{\mathrm{KL}}(P_{\theta_{\mathrm{un}}}|_{S,Z}P_{S,Z} \| \mu^n_{\backslash R} \otimes P_{\theta_{\mathrm{un}}})$
= $I(S, Z; \theta_{\mathrm{un}}),$

where $\mu_{\backslash R}$ denotes the marginal distribution of μ over S and Z (i.e., excluding R).

Hence, we have:

gen
$$(\mu, \mathcal{A}_{un}) \leq \sqrt{\frac{2\sigma^2}{n} \cdot I(S, Z; \theta_{un})}.$$

Note that we use $I(S, Z; \theta_{un})$ instead of $I(S, R, Z; \theta_{un})$, which is the key difference from the typical setup in Xu & Raginsky (2017), where the learned model is assumed to depend on all the data.

989 We can then arrive at the result for training without rationale by noting:

 $I(\theta_{un}; S, Z) = I(\theta_{un}; Z) + I(\theta_{un}; S|Z).$

992 For training with rationale, we have:

$$\begin{split} I(\theta_{\mathrm{ra}}; S, R, Z) &= I(\theta_{\mathrm{ra}}; Z) + I(\theta_{\mathrm{ra}}; R|Z) + I(\theta_{\mathrm{ra}}; S|Z, R) \\ &\leq I(\theta_{\mathrm{ra}}; Z) + H(R|Z) + I(\theta_{\mathrm{ra}}; S|Z, R) \\ &\leq I(\theta_{\mathrm{ra}}; Z) + \eta_1 + \delta, \end{split}$$

where the first inequality is due to $I(\theta_{ra}; R|Z) \le H(R|Z)$, and the second inequality follows from conditions 1 and 2. We can now apply (Xu & Raginsky, 2017, Thm. 1) to yield the result.

A.3 SUPPORTING LEMMAS

Lemma 1 Let \hat{Z} be the estimate of Z based on θ . Let $P_e = P(Z \neq \hat{Z} | \hat{Z}, \theta)$ be the probability of error in predicting Z using θ . Suppose $P_e \leq \epsilon$. Then, we have that:

$$P_e \ge \frac{H(Z) - I(R;\theta) - I(Z;\theta|R) - H(\epsilon)}{\log|Z|}.$$
(9)

Proof: First, let's define an indicator variable E for the error event: $E = \begin{cases} 0, & \text{if } \hat{Z} = Z \\ 1, & \text{if } \hat{Z} \neq Z \end{cases}$. We have

$$\begin{split} H(Z|\theta) &= H(Z|\hat{Z},\theta) = H(Z,E|\hat{Z},\theta) \\ &= H(E|\hat{Z},\theta) + H(Z|E,\hat{Z},\theta) \\ &= H(P_e) + P_e H(Z|E=1,\hat{Z},\theta) + (1-P_e)H(Z|E=0,\hat{Z},\theta) \\ &\leq H(\epsilon) + P_e \log |Z|. \end{split}$$

1014 Hence, we have

$$P_e \ge \frac{H(Z|\theta) - H(\epsilon)}{\log|Z|} = \frac{H(Z) - I(Z;\theta) - H(\epsilon)}{\log|Z|}.$$
(10)

 $\theta|Z)$

1018 This part of the proof follows from Fano's inequality.

Now, since
$$I(Z; \theta) + I(R; \theta|Z) = I(Z, R; \theta)$$
, we have
 $I(Z; \theta) = I(Z, R; \theta) - I(R; \theta|Z)$
 $I(Z; \theta) = I(R; \theta) + I(Z; \theta|R) - I(R; \theta|Z)$

1023
$$\leq I(R;\theta) + I(Z;\theta|R)$$

1024 Plugging in equation 10 obtains the result.

Note: $I(Z; \theta | R) = 0$ if θ is trained only on R, i.e., Z and θ are conditionally independent given R.

Lemma 2 Let \hat{Z} be the estimate of Z based on θ . Let $P_e = P(Z \neq \hat{Z} | \hat{Z}, \theta)$ be the probability of error in predicting Z using θ . Assume that $P_e \leq 0.5$. Then, we have that:

$$P_e \le \frac{H(Z) - I(R;\theta) + H(R|Z)}{\log 2}.$$
(11)

that

Proof: By definition of *E* from Lemma 1, we have

1032
1033
1034
1034
1035
1036
1036
1037
1037
1038
Hence, we have
$$P_e \leq H(Z|\theta)/\log 2$$
. Since $H(Z|\theta) = H(Z) - I(Z;\theta)$, and
1040
 $I(Z;\theta) = I(Z,R;\theta) - I(R;\theta|Z)$

1041
1042
1042
1043
1044
$$= I(R;\theta) + I(Z;\theta|R) - I(R;\theta|Z)$$

$$\ge I(R;\theta) - H(R|Z) + H(R|\theta,Z)$$

$$\ge I(R;\theta) - H(R|Z),$$

the bound follows.

Partial justification of Condition 2: Consider the Markov chain $Z \to R \to \theta$ with additional arrows of $Z \to \theta$ and $R \to \theta$, where Z is the preference, R is the rationale, and θ is a trained model used to predict Z. From both Lemma 1 and 2, we see that as long as θ captures the information of R, i.e., $I(R;\theta)$ is large, and R does not contain excessive irrelevant information other than Z, i.e., H(R|Z) is small, the prediction error of Z from model θ can be well-controlled. Specifically, based on the lower and upper bound analysis, we can conclude that the probability of error P_e decreases with increasing $I(R;\theta)$ or decreasing H(R|Z). This implies that the model does not need to capture additional information from S to achieve high prediction accuracy for Z, i.e., $I(\theta_{ra}; S|Z, R)$ is small. In other words, the incorporation of R in the training of θ_{ra} guides the model to easily predict Z without resorting to finding potentially irrelevant information from S. A full justification of the condition hinges on a detailed analysis of the specific algorithm and is beyond the scope of this study.

A.4 PROOF OF THEOREM 1 AND DERIVATION OF REGIMES

Recall the definition of $\alpha, \beta, \epsilon, f(\cdot)$ in Sec. A.1. Now, we derive the relationship between the mutual information I(Z; R|S) and these parameters:

$$P(Z = 1|S, R = 1) = \frac{P(R = 1|Z = 1, S)P(Z = 1|S)}{P(R = 1|S)}$$
$$= \frac{P(R = 1|Z = 1, S)P(Z = 1|S)}{\sum_{z \in \{0,1\}} P(R = 1|Z = z, S)P(Z = z|S)}$$

$$= \frac{\beta(f(S) + \epsilon)}{\beta(f(S) + \epsilon) + \alpha(1 - f(S) - \epsilon)}$$

Similarly, we have

$$P(Z = 1|S, R = 0) = \frac{\alpha(f(S) + \epsilon)}{\alpha(f(S) + \epsilon) + \beta(1 - f(S) - \epsilon)}$$

These equations show that the probability of Z = 1 given the query X, the preferred response Y_1 , the dispreferred response Y_2 , and the rationale R depends on both the informativeness of the rationale, through α and β , and the informativeness of the query and responses, through f(S). Using the above equations, we get the conditional entropies as follows:

1076
1077
$$H(Z|S, R = 1) = H\left(\frac{\beta(f(S) + \epsilon)}{\beta(f(S) + \epsilon) + \alpha(1 - f(S) - \epsilon)}\right),$$

1078
1079
$$H(Z|S, R = 0) = H\left(\frac{\alpha(f(S) + \epsilon)}{(\epsilon(S) - \epsilon)(\epsilon(S) - \epsilon)$$

$$H(Z|S, R = 0) = H\left(\frac{\alpha(f(S) + \epsilon)}{\alpha(f(S) + \epsilon) + \beta(1 - f(S) - \epsilon)}\right)$$

and substituting to the mutual information equation, we get:

$$I(Z; R|S) = H(Z|S) - \sum_{r=\{0,1\}} P(R=r|S)H(Z|S, R=r).$$

1084 Then, we compute each probabilities P(R|S) as follows:

$$P(R=1|S) = \sum_{z \in \{0,1\}} P(R=1|Z=z,S) P(Z=z|S)$$

$$= \beta(f(S) + \epsilon) + \alpha(1 - f(S) - \epsilon),$$

$$P(R = 0|S) = 1 - P(R = 1|S)$$

1082 1083

1086 1087 1088

1089

1092

$$= 1 - \left(\beta(f(S) + \epsilon) + \alpha(1 - f(S) - \epsilon)\right)$$

and substitute them back into the conditional mutual information term, where we define p = f(S):

$$\begin{split} I(Z;R|S) &= H(p+\epsilon) - (\beta(p+\epsilon) + \alpha(1-p-\epsilon))H\left(\frac{\beta(p+\epsilon)}{\beta(p+\epsilon) + \alpha(1-p-\epsilon)}\right) \\ &- (1 - (\beta(p+\epsilon) + \alpha(1-p-\epsilon)))H\left(\frac{\alpha(p+\epsilon)}{\alpha(p+\epsilon) + \beta(1-p-\epsilon)}\right). \end{split}$$

1099 To study the influence of the parameters α, β, ϵ on the mutual information, we consider the following 1100 edge cases:

1101 1102 Regime 1: Highly informative rationale $\beta \approx 1$ and low noise $\alpha \approx 0 \implies$ Rationale is a sufficient 1103 statistics.

In this case, the conditional probabilities are simplified to

1105
1106
1107
1108
1109

$$P(Z = 1|S, R = 1) \approx \frac{f(S) + \epsilon}{f(S) + \epsilon} = 1,$$

$$P(Z = 1|S, R = 0) \approx \frac{0}{1 - f(S) - \epsilon} = 0.$$

¹¹¹⁰ Thus, the mutual information becomes as follows:

$$\begin{array}{ll} \mbox{1111} \\ \mbox{1112} \\ \mbox{1113} \\ \end{array} & I(Z;R|S) \approx H(f(S)+\epsilon) - P(R=1|S)H(1) \\ -P(R=0|S)H(0) = H(f(S)+\epsilon). \end{array}$$

In this regime, the conditional mutual information is solely determined by the entropy of the preference prediction, $H(f(S) + \epsilon)$. We notice that the entropy function is concave and reaches the max value at the 0.5 mark.

1118 Regime 2: Uninformative rationale $\beta \approx 0.5$ and high noise $\alpha \approx 0.5$.

1121 For this case, the conditional probabilities become:

$$P(Z = 1|S, R = 1) \approx \frac{f(S) + \epsilon}{f(S) + \epsilon + 1 - f(S) - \epsilon} = f(S) + \epsilon = P(Z = 1|S, R = 0)$$

and the mutual information equals to:

$$I(Z; R|S) \approx H(f(S) + \epsilon) - H(f(S) + \epsilon) = 0,$$

which shows that the rationales provides 0 information about the preference given the prompt and responses.

1130 Regime 3: Moderately informative rationale $\beta = 0.5 + \gamma$ and lower noise $\alpha = 0.5 - \gamma$.

1132

1120

1122 1123 1124

1126

Given the assumption of $\beta = 0.5 + \gamma$ and $\alpha = 0.5 - \gamma$, where $0 \le \gamma \le 0.5$ and γ denotes the level of informativeness of the rationale, we substitute this into the conditional mutual information term.

We first substitute into the conditional probabilities and get: $P(Z = 1|S, R = 1) = \frac{(0.5 + \gamma)(f(S) + \epsilon)}{(0.5 + \gamma)(f(S) + \epsilon) + (0.5 - \gamma)(1 - f(S) - \epsilon)},$ $P(Z = 1|S, R = 0) = \frac{(0.5 - \gamma)(f(S) + \epsilon)}{(0.5 - \gamma)(f(S) + \epsilon) + (0.5 + \gamma)(1 - f(S) - \epsilon)}.$ Then, we compute the following probabilities: $P(R=1|S) = \sum_{z \in \{0,1\}} P(R=1|Z=z,S) P(Z=z|S)$ $= 0.5 + 2\gamma (f(S) + \epsilon - 0.5),$ $P(R = 0|S) = 0.5 - 2\gamma(f(S) + \epsilon - 0.5).$ Now, we can compute the conditional mutual information term: $I(Z;R|S) = H(f(S) + \epsilon)$ $-(0.5+2\gamma(f(S)+\epsilon-0.5)) \cdot H\left(\frac{(0.5+\gamma)(f(S)+\epsilon)}{(0.5+\gamma)(f(S)+\epsilon)+(0.5-\gamma)(1-f(S)-\epsilon)}\right)$ $-(0.5 - 2\gamma(f(S) + \epsilon - 0.5)) \cdot H\left(\frac{(0.5 - \gamma)(f(S) + \epsilon) + (0.5 - \gamma)(1 - f(S) - \epsilon)}{(0.5 - \gamma)(f(S) + \epsilon) + (0.5 + \gamma)(1 - f(S) - \epsilon)}\right).$ (12)I(Z;R,S) **Mutual Information** 0.2-0.2 0.4 γ Gamma y Figure 5: The plot of Equation 12 showing the relation between mutual information and gamma γ for a fixed f(S). We can now analyze the behavior of the mutual information as a function of γ : When the rationale is uninformative $\gamma = 0$, then the mutual information becomes 0, I(Z; R|S) = 0,

when the rationale is uninformative $\gamma = 0$, then the mutual information becomes 0, I(Z; R|S) = 0, which is consistent with previous cases, in which uninformative rationales provide no additional information about the preference Z as demonstrated in Figure 5.

As rationale becomes more informative about the preference by increasing γ , we observe that mutual information also increases displayed in Figure 5.

1187 Consider the case that the true probability $P(Z = 1|S) = f(S) + \epsilon > 0.5$, so the preference Z is most likely to be 1.



- For the second entropy term, with an increase of γ , the weight of the second entropy term decreases, and the entropy decreases itself (see Figure 7).
- 1239 1240
- The net effect on the mutual information depends on the relative magnitudes of these changes.However, we can argue that the decrease in the entropy terms dominates the change in their weights

due to the entropy function changing more rapidly near the extremes (i.e., when the distribution is close to being deterministic) compared to the middle range.

Thus, with an increase in γ , the overall contribution of the entropy terms to the mutual information decreases, causing an increase in I(Z; R|S), which indicates that as the rationale becomes more informative about the preference, the mutual information increases.

1248

1250

1252

1262

1263

1249 B ADDITIONAL EXPERIMENTAL RESULTS

1251 B.1 EXPERIMENTAL DETAILS

For DPO-based methods, we fine-tune the base model by supervised fine-tuning (SFT) with the chosen responses from the preference dataset for a single epoch. For ORPO, which avoids the reference model, we skip this SFT step. For models trained on RDPO and results reported in Section 4, we use $\gamma = 2.0$, and for RORPO, we use $\gamma = 10.0$. Similar to baseline methods, we train RDPO and RORPO for 1 epoch. We perform ablation studies on the hyperparameter γ and the number of epochs in the following sections.

For our winrate scores, we report the mean winrates after querying the evaluator 3 times. To reduce the evaluator's order bias, we have additionally shuffled the order of responses. We note that the winrate error bars are within < 3% for 512 samples.

B.2 ABLATION STUDY ON MODELS AND HYPERPARAMETERS

1264						1.55	
1265					M1S	tral-7/E	3-v0.1
1266	General	62	61	63	61	62	60
1267	Detailed	64	66	63	61	60	62
1268	γ	1.0	1.5	2.0	2.5	3.0	10.0
1269	/	110	110		2.10	2.0	1010
1270				Mistra	1-7B-v	0.2-In	struct
1271	C	55	()	57	55	57	50
1272	General	33 56	62	57	33	57	50
1273	Detailed	56	56	57	60	57	59
1274	γ	1.0	1.5	2.0	2.5	3.0	10.0
1275					_		-
1276					Zepl	hyr-7B	B-Beta
1277	General	58	55	55	57	55	57
1278	Detailed	48	49	51	51	50	51
1279	γ	1.0	1.5	2.0	2.5	3.0	10.0
1000	/					- • •	

1281Table 5: The impact of different values of hyperparameter γ on the winrate of the RDPO model1282against the DPO model. The results on various models: Mistral-7B-v0.1 (Top), Mistral-7B-Instruct-1283v0.2 (Middle), and Zephyr-7B-Beta (Bottom).

1284

Here, we investigate the impact of the hyperparameter γ , ranging from 1.0 to 10.0, on the performance of the model trained with RDPO loss. We provide the winrate scores against the DPO model on the Orca dataset. As we see in Table 5, models trained on either general or detailed rationales can still achieve a stronger winrate against the DPO model, except for the case of *Zephyr-7B-Beta* model, which achieves a draw with the DPO model. For this model, the better quality rationales are important for effective preference learning, as the general rationales can still improve the performance.

- 1291
- 1292

1293 B.2.1 RATIONALE-ONLY OPTIMIZATION

1295 In contrast to the ablation study on the rationale hyperparameter, we conduct an exploratory study on an extreme case where the rationale alone or the preference learning alone drives preference

1296		RDPO (Preference + Rationale)	DPO (Preference-Only)	Rationale-Only
1297	General	64.5	59.1	61.8
1290	Detailed	64.4	59.1	61.3

1300 1301

Table 6: The impact of different components of RDPO by measuring the winrate of the target model against the SFT model. The results on Mistral-7B-v0.2-Instruct model using the Orca dataset. 1302

1303

1304 optimization. To address this, we conducted a series of experiments to isolate the impact of each 1305 component and evaluate their combined effect. Specifically, we investigated an extreme case where 1306 the rationale loss alone drives preference optimization, with the DPO alignment loss set to zero. This approach was based on the hypothesis that rationales inherently encode preferences by combining 1307 preference-response pairs, the preferences themselves, and the associated reasoning processes, thereby 1308 providing a rich and effective training signal. 1309

1310 For these experiments, we fine-tuned Mistral-7B-Instruct-v0.2 on the Orca dataset across three 1311 settings: RDPO (combining DPO and rationale loss), DPO (excluding rationale loss), and Rationale-1312 Only (excluding DPO loss). The results, as shown in the Table 6, reveal that rationales alone can 1313 substantially improve model performance, achieving a high win rate of over 61While DPO also demonstrated a majority win rate against the SFT baseline, training with both rationale and preference 1314 losses (RDPO) consistently achieved the highest win rate (64.5%) across both general and detailed 1315 settings. This highlights the benefit of integrating rationales into the preference objective, effectively 1316 leveraging the strengths of both losses to produce superior performance. To further investigate 1317 how rationales enhance DPO preference learning, we examined the reward margin metrics. As 1318 shown in the Table 7, RDPO not only achieved higher reward margins between chosen and rejected 1319 responses but also demonstrated faster convergence compared to DPO. This can be explained through 1320 the following: while DPO explicitly aims to maximize reward margins, the inclusion of rationales 1321 provides an implicit quality signal, offering explanations for the differences between chosen and 1322 rejected responses. This signal reinforces the model's ability to improve reward margins by guiding it 1323 toward more informed preferences.

24													
25	Training Points	0	1000	2000	3000	4000	5000	6000	7000	8000	9000	10000	11000
,)	DPO	0.00	0.05	0.19	0.32	0.42	0.49	0.54	0.58	0.62	0.63	0.65	0.66
	RDPO	0.00	0.10	0.25	0.46	0.67	0.76	0.83	0.85	0.86	0.87	0.89	0.91

1328 Table 7: The eval reward margin comparison for different training losses, DPO and RDPO. The 1329 Mistral-7B-Instruct-v0.2 is trained on the Orca dataset. 1330

1331 These findings underscore the complementary nature of the rationale SFT loss and the pairwise 1332 alignment loss. While DPO explicitly optimizes reward margins, the rationale prediction loss 1333 provides supplementary supervision, enabling the model to learn the reasoning underlying response 1334 preferences. This integration not only strengthens the selection process but also accelerates training 1335 convergence. By combining these two approaches, RDPO amplifies their individual strengths, 1336 resulting in more efficient and effective preference learning.

1337 1338

1339

B.3 ABLATION ON THE NUMBER OF EPOCHS

1340		General	Detailed
1341	Epoch 1	76	74
1342	Epoch 2	75	71
1343	Epoch 3	74	74

Table 8: The analysis of the number of epochs on the performance. Winrate of the RDPO model 1345 trained on the general rationales (left) and detailed rationales (right) against the DPO model, respec-1346 tively. 1347

1348

1344

In the main paper, we trained the models with the RDPO loss on a single epoch similar to 1349 DPO (Rafailov et al., 2024). Here, we study the impact of training the models with the rationales on more epochs. As we observe in Table 8, training with more epochs does not improve the winrate of the RDPO model against the DPO model. The reason could be that the model has already learned the preference well after the first epoch, or it could be that the quality of rationales could be further improved to increase the efficiency of the rationale-based preference learning algorithms.

B.4 EVALUATION OF THE IMPACT OF RATIONALES WITH ALPACAEVAL 2.0

Here, we study the flexibility of our method by extending different preference learning methods, such as DPO (Rafailov et al., 2024) and ORPO (Hong et al., 2024), with rationales. In this experiment, we evaluate the performance of the trained models on the automatic instruction-following AlpacaEval 2.0 benchmark (Li et al., 2023b) with GPT-4-turbo as a judge and report the raw winrate, the length-controlled (LC) winrate (Dubois et al., 2024), which is robust against the verbosity bias that the raw winrate inherently entails, and the average response length. We train the instruction-tuned models, Mistral-7B-Instruct-v0.2 and Llama-3-8B-Instruct, on the Intel-DPO-Pairs dataset.



Figure 8: The performance comparison of the original model, DPO trained model, and DPO with
rationales (RDPO) trained model on the AlpacaEval 2.0 benchmark. The **bolded** numbers denote the
average response length of each models.

We observe in Figure 8 that DPO trained model improves the winrates on both models compared to the original model. Additionally, RDPO models further increase the win rates. We note in Figure 9 that in the case of the Mistral model, the winrate decrease with ORPO preference training, which might be due to the lack of the reliance on the reference model, the behavior which is also observed in Ethayarajh et al.. Furthermore, after adding rationales, the (LC) winrate not only increases but also surpasses of the original model. These results show the helpfulness of adding rationales into preference learning.

Interestingly, we also note that rationale based models increase the winrates while their average response lengths are decreased compared to average response lengths of the original model, which is not a similar observation as seen in some current methods, such as SPPO (Wu et al., 2024) or SimPO (Meng et al., 2024).

1400 B.5 COST ANALYSIS

In this section, we analyze the cost breakdown to assist project owners in evaluating trade-offs.
 Specifically, we present a cost-benefit analysis of the approach. The table below outlines the cost of using the API to generate rationales for a given number of annotations. It also highlights the RDPO



1458 **B.6** RATIONALE GENERATION 1459

1460 1461

1462

1463

1480

1481

1491

1493

Here, we provide the prompts to generate detailed and general rationales as mentioned in Section 4. 1464 In Prompt 1, we ask the model to generate general rationales without going into details and provide 1465 high-level reasoning. In Prompt 2, we ask the model to provide a specific rationale and go as much 1466 into detail as possible. These are our examples of prompts and we call for further studies to improve 1467 the prompts for higher quality rationales. 1468

Given the prompt: 1469 <prompt> 1470 Which of the following assistant's responses is preferred and 1471 strictly follows the prompt question? Why this assistant's 1472 response: 1473 <chosen> 1474 is preferred over the response provided below: 1475 <rejected> 1476 On a high and general level, why the response above is 1477 preferred over the response below? Provide a general, 1478 high-level explanation for your reasoning without going into the response's details. 1479

Prompt 1: Prompt to the LLM to retrieve a general rationale.

1482 Given the prompt: <prompt> 1483 Which of the following assistant's responses is preferred and 1484 strictly follows the prompt question? Why this assistant's 1485 response: 1486 <chosen> 1487 is preferred over the response provided below: 1488 <rejected> 1489 Why the response above is preferred over the response below? 1490 Provide a specific, detailed explanation for your reasoning pointing directly at the response. 1492

Prompt 2: Prompt to the LLM to retrieve a detailed rationale.

B.7 COMPARISON OF RATIONALES

1504

1506

1507

We provide examples of the general and detailed rationales for the Orca dataset by the Mistral-7B-Instruct-v0.2, and compare them in Examples 1, 2, and 3. As we observe, the general and detailed 1509 rationales follow our intent. Additionally, we observe that a part of the general rationale is also 1510 included in the detailed rationale. It will be important to study, which type of rationale is more helpful 1511 for the preference training.

1512	<general rationale=""></general>
1513	The first response strictly adheres to the prompt by providing
1514	the RDF triplets in the required format, while the second
1515	response includes an explanation that goes beyond the required
1516	format and may not be necessary for the task at hand. Thus,
1517	the first response is preferred as it focuses solely on
1518	providing the required output in the correct format.
1519	
1520	<detailed rationale=""></detailed>
1521	The assistant's response is preferred because it follows the
1522	standard format for RDF triplets, where the subject comes
1523	first, followed by the predicate, and then the object. In the
1524	given input sentence, 'AFC Ajax (amateurs)' and 'Ajax Youth
1525	Academy are the subjects, has ground and prays at are the
1526	predicates, and sportpark be rockomst is the object that is
1527	structured as [subjects predicate object] with each triplet
1522	representing a distinct relationship. The assistant's response
1520	correctly reflects this structure and semantics making it a
1529	preferred option
1030	protottoa operon.
1531	Example 1: An example $\#1$ of general and detailed rationales.
1532	Conoral Pationalo
1533	The preferred response strictly follows the prompt question by
1534	calculating the total amount naid per episode for both major
1535	and minor characters based on the given information. It uses
1536	the provided information to calculate the salaries of the major
1537	and minor characters and then adds them together to find the
1538	total amount paid per episode. The response below includes
1539	unnecessary details and makes the response longer and less
1540	efficient.
1541	
1542	<detailed rationale=""></detailed>
1543	The assistant's response above is clear and accurate. They
1544	first calculate the total amount paid to the minor characters
1545	per episode, which is \$60,000. Then, they calculate the
1546	total amount paid to the major characters per episode, which
1547	is \$225,000. Finally, they add the salaries of the major
1548	and minor characters to find the total amount James pays per
1549	episode, which is \$285,000." This explanation directly points
1550	Lo the steps taken in the assistant's response and explains
1551	to make the explanation easy to understand
1552	to make the exprendition easy to understally.
1553	Example 2. An example $\#2$ of general and detailed rationales
1557	Example 2. All example $\#2$ of general and detailed rationales.
1555	
1555	
1550	
1007	
1000	
1559	
1560	
1561	
1562	
1563	
1564	

1000	<pre><general rationale=""></general></pre>
1567	The response above is preferred because it directly answers the
1568	question by specifically identifying the kidneys as the organs
1569	that secrete hormones to help maintain homeostasis, while
1570	the response below goes into excessive detail and provides
1571	additional information that is not directly related to the
1572	question.
1573	
1574	<detailed rationale=""></detailed>
1575	The assistant's response is preferred because it provides
1576	a clear and detailed explanation of the hormones secreted
1577	by the kidneys and their roles in maintaining homeostasis.
1578	The response also uses the terminology from the paragraph,
1570	such as 'erythropoietin' and 'renin,' which demonstrates a
1520	thorough understanding of the topic. In contrast, the less
1500	preferred response only mentions that the kidneys secrete
1501	their functions. Therefore, the preferred response is more
1582	informative and accurate making it a better answer for the
1583	given question
1584	given quescion.
1585 1587 1588 1589 1590	Example 3: An example #3 of general and detailed rationales. winrate using HALOs repository Ethayarajh et al. (2023) B.8 EVALUATING RESPONSES WITH LLM AS A JUDGE
1591	Here is a conversation between <user> and <assistant>.</assistant></user>
1592	<history></history>
1593	<assistant> can now choose from different responses.</assistant>
1594	Response 1: <response1></response1>
1595	END OF RESPONSE 1
1596	Response 2: <response2></response2>
1597	END OF RESPONSE 2
1598	Between these two responses, which response is more helpful,
1599	narmiess, and concise? Pick one.
1000	

Prompt 3: Prompt to the LLM (e.g., GPT-40) to retrieve the preference over two responses. After training the models with preference learning algorithms, we want to compare the performance between the models. Thus, a common practice is to generate the responses from respective models and compare the quality between them. As a human evaluator might be expensive to acquire, we usually use a capable language model to judge the responses. In our case, we leverage the GPT-40 model to judge and provide winrate scores. We provide the prompt to determine a preferred response between generations in Prompt 3, which is implemented in the HALOs code implementation (Ethayarajh et al., 2023) and shuffles responses to avoid the ordering bias.

1608 1609

1566

1610 B.9 DPO AND RDPO GENERATED RESPONSES

1611 We provide some examples and compare responses generated by the models trained with DPO and 1612 RPDO losses, respectively. In Example 5, we observe the case when DPO demonstrates the verbosity 1613 and generates a lengthy output compared to the response by the RDPO model. The length bias in 1614 DPO is also demonstrated in Example 4, which generates not only a longer output than RDPO but 1615 also outputs a personal opinion on a political stance. Example 6 asks for generating a tweet for a 1616 single sentiment. While the RDPO-trained model correctly recognizes the sentiment and creates a 1617 tweet for the selected sentiment, the DPO-trained model produces tweets for all sentiments regardless of the prompt. Despite the responses generated by RDPO being shorter on average than the ones by 1618 DPO, when prompted to generate an article, the RDPO response can be as lengthy as the DPO one as 1619 shown in Example 7.

1620 <Prompt> 1621 What is your opinion on Britain exiting the EU? 1622 What is your opinion on Britain leaving the EU and Brexit? 1623 Would you say that these questions are the same? 1624 <DPO Response> 1625 These questions are related but not exactly the same. Both 1626 ask for an opinion on Britain's decision to leave the European 1627 Union (EU), which is known as Brexit. However, the first 1628 question asks specifically about Britain exiting the EU, while 1629 the second question expands to include the broader concept of 1630 Brexit, which includes not only leaving the EU, but also the 1631 negotiations and potential outcomes of that decision. 1632 1633 Britain exiting the EU refers to the actual act of leaving 1634 the EU, which includes leaving the EU's political and economic structures. Brexit, on the other hand, encompasses the entire 1635 process, including the negotiations, the potential outcomes, 1636 and the impact on British politics and society. 1637 1638 So, if someone is asking for your opinion on Britain exiting 1639 the EU, they are asking for your thoughts on the actual act 1640 of leaving the EU. If they ask for your opinion on Brexit, 1641 they might be asking for your thoughts on the entire process, 1642 including the negotiation period and the potential outcomes of 1643 leaving the EU. 1644 1645 My personal opinion is that Brexit has created uncertainty and potential risks for both the EU and the UK. However, it 1646 also gives the UK the opportunity to set its own trade and 1647 regulatory policies, which could potentially lead to new 1648 opportunities and economic growth. 1649 1650 <RDPO Response> 1651 Yes, I would consider the questions to be similar, but not 1652 exactly the same. The first question asks for your opinion 1653 on Britain's departure from the EU, while the second question 1654 specifically refers to Brexit, which is the name given to the 1655 process of Britain's departure from the EU. So, while the 1656 second question is a subset of the first, it adds a layer of specificity by referring to the name given to the process. 1657 1658 Example 4: A comparison of generated responses from DPO and RDPO models. 1659 1660 1661 1662 1663 1664 1665 1666 1667 1668 1669 1670

- 1671
- 1672
- 1673

1674 <Prompt> 1675 I ' ve already decided that I want to live in Finland for a 1676 semester . What then ? What if - I know , I ' m getting ahead 1677 of my own dream - world - he and I are together for so long ... we fall in love ? What then ? 1678 Answer the following question: Where is the narrator likely to 1679 be in a few months ? 1680 1681 <DPO Response> 1682 In a few months, the narrator is likely to be living in Finland 1683 for a semester as they have already decided on that. 1684 As for their second question, if they fall in love with someone 1685 during their time in Finland, the outcome could vary. The 1686 narrator could choose to continue their relationship despite 1687 the distance if they are committed to making it work. They 1688 could also choose to pursue the relationship further if both parties are interested. 1689 However, it's important to note that the narrator's 1690 hypothetical situation might not always align with reality. Τf 1691 they do fall in love and pursue a long-distance relationship, 1692 it could require more effort, patience, and communication 1693 to maintain the relationship. Ultimately, it depends on the 1694 individuals involved and their level of commitment to making it 1695 work. 1696 In any case, their current plan is to go to Finland for a 1697 semester, so that is where they are likely to be in a few 1698 months. 1699 Regarding the world reference in the second question, it's not clear how it relates to the overall context of the scenario. 1700 If it refers to a fictional world, it's likely that the 1701 question is asking about the potential outcomes of an imaginary 1702 situation. 1703 However, if it refers to the real world, it's worth noting that 1704 long-distance relationships can be challenging, but they are 1705 also possible with dedication, trust, and open communication. 1706 If the narrator decides to pursue a relationship with someone, 1707 they can consider using technology to stay connected, such as 1708 video calls, messaging, and email. They could also set aside 1709 time for regular updates on each other's lives, and plan visits 1710 to see each other as often as possible. Ultimately, the key to maintaining a successful long-distance 1711 relationship is to approach it with a positive attitude, 1712 patience, and an understanding that challenges will arise. 1713 By staying committed to each other and working together 1714 to overcome any obstacles, the narrator and their partner 1715 could potentially build a strong and lasting connection. So. 1716 assuming the narrator's plan to go to Finland for a semester 1717 remains unchanged, they are likely to be in Finland in a few 1718 months. Whether or not they fall in love, the outcome depends 1719 on the individuals involved and their willingness to make the 1720 relationship work despite the challenges. 1721 1722 Response continues... 1723 1724 <RDPO Response> 1725 The narrator is likely to be living in Finland for a semester. However, if they fall in love with someone, their plans might 1726 change depending on their feelings and circumstances. 1727

1728	Example 5: A comparison of generated responses from DPO and RDPO models.
1729	<prompt></prompt>
1731	Sentiment possibilities. Choices:
1732	(A). negative;
1733	(B). positive;. Generate a tweet that has the following
1734	Sentiment: (A).
1735	<dpo response=""></dpo>
1736	(A). negative: "I just received a faulty product from the
1737	company, and it's really frustrating! #CustomerServiceFail"
1738	(B). positive: "Had an amazing time at the concert last night!
1739	The energy was contagious and I'm already counting down the
1740	days for the next one!
1741	(RDPO Response)
1742	Here's a tweet that conveys a negative sentiment:
1743	"Feeling down today. The rain outside matches my mood.
1744	
1745	Example 6: A comparison of generated responses from DPO and RDPO models.
1746	
1747	
1748	
1749	
1751	
1752	
1752	
1754	
1755	
1756	
1757	
1758	
1759	
1760	
1761	
1762	
1763	
1764	
1765	
1767	
1760	
1760	
1770	
1771	
1772	
1773	
1774	
1775	
1776	
1777	
1778	
1779	
1780	
1781	

1782 <Prompt> 1783 Write an article based on this summary: 1784 Jack Tector and Nathan Smith have received their first Ireland 1785 call-up for this month's ICC Intercontinental Cup game against the Netherlands. 1786 1787 <DPO Response> 1788 Title: Fresh Blood in the Irish Squad: Jack Tector and Nathan 1789 Smith Await Their Debut against the Netherlands 1790 1791 Ireland, a long-standing cricket force in European cricket, is 1792 set to feature two of its exciting young talents, Jack Tector 1793 and Nathan Smith, in this month's ICC Intercontinental Cup 1794 game against the Netherlands. These young guns will be making 1795 their debut on the international stage, providing them with an 1796 unforgettable experience and an opportunity to shine amongst the cricketing world's elite. 1797 1798 Jack Tector, a 23-year-old all-rounder, has been a consistent 1799 performer in both domestic and provincial cricket, earning 1800 him a well-deserved Ireland call-up. With his skills in 1801 both batting and bowling, Tector offers immense potential to 1802 the side and can add valuable depth to the Irish squad. His 1803 aggressive approach at the crease, coupled with his ability to 1804 generate pace and bounce from the ball, make him a formidable 1805 prospect to look out for. 1806 1807 Nathan Smith, 21, also makes his international debut after showcasing his prowess in various domestic competitions. He 1808 is known for his deceptive left-arm spin, which has troubled 1809 numerous batsmen in the domestic circuit. 1810 1811 <RDPO Response> 1812 Ireland Calls Up Jack Tector and Nathan Smith for Upcoming ICC 1813 Intercontinental Cup Game against Netherlands 1814 1815 The Irish cricket team has made some exciting announcements 1816 recently, as two young and promising players, Jack Tector and 1817 Nathan Smith, have been granted their first call-ups to the 1818 national squad for the ICC Intercontinental Cup game against the Netherlands. This upcoming match is set to take place this 1819 month, and both Tector and Smith are eager to make the most of 1820 this incredible opportunity. 1821 1822 Jack Tector, a 21-year-old all-rounder from Wiltshire, England, 1823 has been on the radar of Irish selectors for some time now. He 1824 made his first-class debut for the Northamptonshire Steelbacks 1825 in 2019, and since then, he has been making waves in domestic 1826 cricket competitions, particularly in the 50-over format. 1827 Tector has been consistently impressive with the bat, having 1828 scored a century and several half-centuries in various domestic 1829 games. He is also a handy left-arm spinner, which adds depth to Ireland's bowling attack. 1830 1831 Nathan Smith, on the other hand, is a 23-year-old left-handed 1832 opening batsman from County Antrim, Northern Ireland. 1833 1834 Example 7: A comparison of generated responses from DPO and RDPO models.

1835

1836 B.10 HARDWARE AND SOFTWARE

We conduct our experiments on 3 x Nvidia A100 GPUs. Training $\leq 8B$ -parameter models on preference learning with 12,000 training data points requires $\sim 2 - 3$ GPU hours while generating responses from 512 test samples takes around 10 minutes. We publish our code implementation for generating rationales and training with rationale-enriched preference learning in the anonymous repository for the time of review: https://anonymous.4open.science/r/rationale-3973/.