# **Removing Length Bias in RLHF Is Not Enough**

Anonymous Author(s) Affiliation Address email

### Abstract

| 1  | Reinforcement Learning from Human Feedback (RLHF) has become an essential              |
|----|--|
| 2  | technique for enhancing pretrained large language models (LLMs) to generate            |
| 3  | responses that align with human preferences and societal values. While RLHF has        |
| 4  | shown promise, the training of reward models (RMs) still faces the challenge of        |
| 5  | reward hacking, motivating recent works to prevent RMs from finding shortcuts          |
| 6  | that bypass the intended optimization objectives by identifying simplistic patterns,   |
| 7  | especially response length. Besides the issue of length bias, our work firstly reveal  |
| 8  | that prompt-template bias learned by RMs can also cause reward hacking when            |
| 9  | dealing with marginal samples, resulting in LLMs preferring to generate responses      |
| 10 | in a specific format after RLHF fine-tuning, regardless of the format requested in the |
| 11 | prompt. To this end, we propose a low-cost but effective method, namely Prompt         |
| 12 | Bias Calibration (PBC), to estimate the <i>prompt-template bias</i> term during reward |
| 13 | modeling, which can be utilized to calibrate reward scores in the following RL         |
| 14 | fine-tuning process. Then, we show that our PBC method can be flexibly combined        |
| 15 | with existing algorithms of removing length bias, leading to a further improvement     |
| 16 | in the aspect of enhancing the quality of generated responses. Experiments results     |
| 17 | show that the performance of our PBC method and its extensions have significantly      |
| 18 | surpassed the original implementation of RLHF.   |

### 19 1 Introduction

Reinforcement Learning from Human Feedback (RLHF) has become a critical technique to enable 20 pretrained large language models (LLMs) to follow human instructions, understand human intent, 21 and also generate responses that align with human preferences and societal values [1–4]. Specifically, 22 RLHF usually trains a reward model (RM) to act as the proxy of human preferences, and then 23 utilize online reinforcement learning (RL) algorithms to fine-tune the language models for generating 24 25 responses that can achieve higher expectation rewards, leading to the success of ChatGPT and also 26 many other AI systems [5, 6]. Although the paradigm of RLHF has simplified human data collection, as acquiring human ratings is much easier than collecting demonstrations for supervised fine-tuning 27 (SFT), it still requires huge amount of human-annotated preference pairs to train well-performing 28 RMs in practice, motivating recent researches to seek novel alignment methods to bypass RM 29 training [2-4]. However, the pipeline of original RLHF is still the primary choice of most industrial 30 applications, because well-trained RMs can provide a certain level of generalization ability [7]. 31

Besides the expensive cost of collecting numerous human-annotated preference pairs, another heavily criticized issue of RLHF could be the phenomenon of *reward hacking* [8], where the over-optimized RMs tend to find some shortcuts to bypass its intended optimization objective, through identifying some simple patterns to distinguish between good and bad responses [9]. The most widely studied pattern in *reward hacking* could be the sentence (response) length, and these trained RMs can utilize the preference among human raters for longer responses to achieve *reward hacking*, despite the actual

quality of response does not improve with the increase of response length [10]. Thus, to mitigate

*reward hacking*, recent works has primarily focused on estimating the *length bias* term in the reward scoring process, so that it can be removed in the subsequent RL fine-tuning procedure to further

<sup>41</sup> improve the quality of generated response after RLHF process [11, 12].

42 Besides the issue of *length bias*, in the practice of applying RLHF to industrial products, we have 43 observed that the original implementation of RLHF tends to make LLMs prefer generating responses 44 in a specific format. This observation motivates us to investigate the underlying causes and seek a 45 cost-effective solution to address this issue. The main contributions are summarized as follows:

- We are the first to reveal the existence of *prompt-template bias* in RMs trained with the original preference loss, and theoretically analyze the cause of *prompt-template bias* issue, along with its corresponding potential risks on the entire RLHF process;
- To mitigate the *reward hacking* caused by *prompt-template bias*, we develop a Prompt Bias
   Calibration (PBC) method, which will firstly estimate the *prompt-template bias* term during
   the reward scoring process, and then remove it in the subsequent RL fine-tuing process;
- We show that the developed PBC method can be flexibly combined with existing methods of removing *length bias*, leading to a further improvement in the aspect of enhancing the quality of generated responses;
- Experimental results show that our developed PCB method and its extensions can achieve promising performance improvements compared to the original implementation of RLHF.

### 57 2 Preliminary

Reward models (RMs) have become the dominant tool for aligning the LLM's responses with user
preferences or task-specific requirements [1, 9]. In this section, we will firstly review the training
procedure of reward models in Sec. 2.1, including analyzing the causes of *length bias* and *prompt bias* in existing RMs, and also illustrate how these RMs are used for alignment in Sec. 2.2, especially
RLHF fine-tuning processes.

#### 63 2.1 Reward Model Training

The usual optimization goal of a reward model is to minimize the loss under the Bradley–Terry model [13] on the dataset of pair-wise comparisons of model responses, denoted as  $(x, y^+, y^-) \in \mathcal{D}$  where *x* indicates the input prompt,  $y^+$  and  $y^-$  are the chosen and rejected responses respectively. Then, the objective function can be formulated as

$$\mathcal{L}^{RM}(\theta) = -\mathbb{E}_{(x,y^+,y^-)\sim\mathcal{D}}\left[\log(\sigma(r_\theta(x,y^+) - r_\theta(x,y^-)))\right]$$
(1)

where  $r_{\theta}(x, y)$  denotes the reward model that takes the prompt x and response y as input to predict a scalar reward with trainable parameters  $\theta$ ;  $\sigma$  denotes the sigmoid function.

<sup>70</sup> Length Bias: Denote  $r_{\theta^*}(x, y)$  as the "gold standard" reward model [9] with the optimal parameters <sup>71</sup>  $\theta^*$ , it reflects human's intrinsic ranking preferences and can play a role of human rater to provide gold <sup>72</sup> reward signal for each prompt-response pair. However, due to the subjectivity of ranking preferences <sup>73</sup> and flaws in rating criteria, there is a phenomenon where human raters prefer longer responses that <sup>74</sup> appear to be more detailed or better formatted, but their actual quality does not improve [10]. Thus, <sup>75</sup> the "gold standard" reward model for rating preference data can often be biased and thus we can <sup>76</sup> decompose it to disentangle the actual reward from the spurious reward [11], formulated as

$$r_{\theta^*}(x,y) = r_{\theta^*}^Q(x,y) + r_{\theta^*}^L(x,y),$$
(2)

where  $r_{\theta^*}^Q(x, y)$  is the actual reward gains brought by improving the quality of response y;  $r_{\theta^*}^L(x, y)$ is the spurious reward gains of increasing response length, whose patterns are much easier to identify.

Thus, with *length bias* in the "gold standard"  $r_{\theta^*}(x, y)$ , during the training of reward model,  $r_{\theta}(x, y)$ can easily find shortcuts to bypass its intended optimization objective, through identifying simple patterns, such as sentence (response) length, to distinguish between good and bad responses, leading to the phenomenon of "reward hacking" caused by *length bias* [10]. Without increasing the cost of rating higher quality preference data, it becomes increasingly important and beneficial to study mitigating the impact of *length bias* in the process of reward modeling.



Figure 1: Comparison of the RM training process using the original preference loss and our developed PBC method respectively, where the latter employs  $u_c(x)$  to approximate the *prompt-template bias*, providing unbiased reward scores with lower variance for the subsequent RL fine-tuning.

- Prompt Bias: the prompt bias in reward modeling derives from the underdetermination of Bardley-
- <sup>86</sup> Terry model [13]. For any reward model  $r_{\theta'}(x, y)$  learned from the preference loss defined in Eq. (1),
- whose target is optimized to approximate the "gold standard"  $r_{\theta^*}(x, y)$ , there always exists an
- equivalent reward model  $r_{\theta}(x, y)$  that satisfies

$$r_{\theta}(x,y) := r_{\theta'}(x,y) + C(x) \tag{3}$$

where C(x) is a prompt-dependent constant referred to as *prompt bias*, leading to the same loss value as  $\mathcal{L}(\theta) = \mathcal{L}(\theta')$ . Due to the fact that there is no constraint on C(x) in the original preference loss as defined in Eq. (1), the issue of *prompt bias* has been criticized in the scenario of reward model ensembles [8], where different reward models tend to choose different values for C(x), making the statistics of the set of reward scores meaningless.

As shown in Fig. 1, it has been widely reported that the *prompt bias* will result in a certain gap in the
 mean values of the set of prompt-response pairs under different prompts. However, in our research,
 we find that this gap is more likely caused by the *prompt-template bias*, as discussed in Section 3.1.

#### 97 2.2 RLHF Fine-tuning

Given the trained reward model  $r_{\theta}(x, y)$  as the proxy of human preferences, Reinforcement Learning from Human Feedback (RLHF) tends to utilize an online reinforcement learning method, typically proximal policy optimization (PPO) [14], trains a policy language model  $\pi_{\phi}^{RL}$  to maximize expected reward, while staying close to its initial policy  $\pi_{\phi}^{SFT}$ , which is finetuned on supervised data (promptresponse pairs). Through measuring the distance from the initial policy with Kullback-Leibler (KL) divergence, the optimization objective of RLHF fine-tuning can be formulated as

$$\mathcal{L}^{RL}(\phi) = \mathbb{E}_{(x,y)\sim\mathcal{D}_{\pi_{\phi}^{RL}}}\left[r_{\theta}(x,y) + \beta \log\left[\pi_{\phi}^{RL}(y|x)/\pi^{SFT}(y|x)\right]\right],\tag{4}$$

where  $\beta$  is the hyper-parameter to control the strength of the KL divergence term.

### 105 3 Method

In this section, we will firstly investigate the cause of *prompt-template bias* and then theoretically analyze its potential risks when dealing with marginal samples during reward modeling, as shown in Sec. 3.1, and then illustrate our low-cost but effective method to estimate the *prompt-template bias* term during RM training in Sec. 3.2, which can be utilized to calibrate reward scores in the following RL fine-tuning process. At last, in Sec. 3.3, we show that our Prompt Bias Calibration (PBC) method can be flexibly combined with recent popular methods of removing *length bias*, leading to a further improvement in the aspect of enhancing the quality of generated responses.

#### 113 3.1 Impact of prompt-template bias on RLHF

In this part, we will first illustrate the cause of prompt-template bias during RM training. Formally, 114 given a set of prompt-response pairs, denoted as  $\mathcal{D}_a = \{x_a, y_a^{(i)}\}_{i=1}^{N_a}$ , with the same user prompt  $x_a$ , 115 e.g. "writing an academic paper on the field of computer science", and  $\{y_a^{(i)}\}_{i=1}^{N_a}$  denoting the set of collected academic papers to satisfy the request of  $x_a$ , the prompt bias term, specifically  $C(x_a)$ , 116 117 learned by RMs is supposed to not affect the preference order within  $\mathcal{D}_a$ , as discussed in Section 2.1. 118 However, in the practice of RM training, the reward score is usually predicted by a LLM that takes 119 the concatenation of the prompt and response as input, making it challenging for RMs to learn a bias 120 term that focuses solely on the prompt x while disregarding variations in the subsequent response y. 121 During the training process to order the pairs within  $\mathcal{D}_a$ , we find that RMs trained with the original 122 preference loss in Eq. (1) are more likely to introduce a joint bias term across the entire sequence of 123 concatenating the prompt and response, formulated as 124

$$r_{\theta}(x_a, y_a) := r_{\theta'}(x_a, y_a) + C(x_a, \overline{y}_a), \quad \overline{y}_a = \frac{1}{N_a} \sum_{i=1}^N y_a^{(i)}, \tag{5}$$

where  $\overline{y}_a$  can be considered the average response of the response set  $\{y_a^{(i)}\}_{i=1}^{N_a}$ , and it will embody the common characteristics found within these collected responses, such as the format of *academic paper*;  $C(x_a, \overline{y}_a)$  denotes the joint bias on the entire sequence of the prompt  $x_a$  associated with the average response  $\overline{y}_a$  in the format of *academic paper*;  $r_\theta(x_a, y_a)$  is still supposed to approximate the "gold standard" provided by  $r_{\theta^*}(x_a, y_a)$ , leading to  $\mathbb{E}_{\mathcal{D}_a} [r_{\theta^*}(x_a, y_a)] \approx \mathbb{E}_{\mathcal{D}_a} [r_{\theta^*}(x_a, y_a)]$ .

Considering the average response  $\overline{y}$  can be treated as a standard template of the response to the 130 prompt x, we define the joint bias  $C(x, \overline{y})$  as prompt-template bias. Then, we highlight the properties 131 of prompt-template bias as follows: 1) the original preference loss in Eq. (1) imposes no constraints 132 on  $C(x, \overline{y})$ , because its value will not influence the outcome of the preference loss and also not affect 133 the preference order within the prompt-response pairs collected for the same prompt x; 2)  $C(x, \overline{y})$ 134 will reduce to the original prompt bias C(x, -) when no common characteristics can be found across 135 all of these collected responses, indicating the diversity of  $\{y^{(i)}\}_{i=1}^N$  is sufficiently high. With these 136 properties in mind, we assume that the prompt-template bias  $C(x, \overline{y})$  can essentially meet most of the 137 properties of the original prompt bias C(x, -) as discussed in Section 2.1. Thus, we suppose  $C(x, \overline{y})$ 138 can be considered as a broader definition of prompt bias in the actual RM training, because it is more 139 likely to be learned by RMs in practice, given the fact that preference pairs are extremely scarce and 140 the diversity of responses collected for the same prompt is often insufficient. 141

After defining prompt-template bias, we will theoretically investigate the impact of introducing 142  $C(x,\overline{y})$  during RM training on the entire RLHF process. Assume that there exist two sets of prompt-143 response pairs, denoted as  $\mathcal{D}_a = \{x_a, y_a^{(i)}\}_{i=1}^{N_a}$  and  $\mathcal{D}_b = \{x_b, y_b^{(i)}\}_{i=1}^{N_b}$ , where  $x_a$  and  $x_b$  indicate different categories of prompts, *e.g.*  $x_a$  requests "writing an **academic paper** on theme **a**" and  $x_b$ 144 145 requests "writing a **brief** on theme **b**", and  $\{y_a^{(i)}\}_{i=1}^{N_a}$  and  $\{y_b^{(i)}\}_{i=1}^{N_b}$  denote the collected responses for answering the prompt  $x_a$  and  $x_b$  respectively. After RM training, due the fact that there is no 146 147 constraint on  $C(x, \overline{y})$  in the preference loss defined in Eq. (1), the discrepancies of prompt biases 148 between these two previously mentioned sets of prompt-response pairs, specifically  $\mathcal{D}_a$  and  $\mathcal{D}_b$ , could 149 be extremely large, e.g.  $C(x_a, \overline{y}_a) >> C(x_b, \overline{y}_b)$ , leading to 150

$$\mathbb{E}_{(x_a, y_a) \sim \mathcal{D}_a} \left[ r_\theta(x_a, y_a) \right] \gg \mathbb{E}_{(x_b, y_b) \sim \mathcal{D}_b} \left[ r_\theta(x_b, y_b) \right] \tag{6}$$

where  $r_{\theta}(x_a, y_a) = r_{\theta'}(x_a, y_a) + C(x_a, \overline{y}_a)$  and  $r_{\theta}(x_b, y_b) = r_{\theta'}(x_b, y_b) + C(x_b, \overline{y}_b)$ . The unbiased reward distributions, modeling the reward scores  $\{r_{\theta'}(x_a, y_a^{(i)})\}_{i=1}^{N_a}$  and  $\{r_{\theta'}(x_b, y_b^{(i)})\}_{i=1}^{N_b}$  respectively, should exhibit similar mean values, *e.g.*  $\mathbb{E}_{\mathcal{D}_a}[r_{\theta'}(x_a, y_a)] \approx \mathbb{E}_{\mathcal{D}_b}[r_{\theta'}(x_b, y_b)]$ , and will make little impact on the comparison of expectation terms in Eq. (6). We highlight that the discrepancies of *prompt bias* terms, specifically the gap between  $C(x_a, \overline{y}_a)$  and  $C(x_b, \overline{y}_b)$ , won't affect preference ordering within categories, but can cause disaster when dealing with some marginal samples, like "*an academic paper on theme b*" denoted as  $y_{ab}$ , or "*a brief on theme a*" denoted as  $y_{ba}$ .

To facilitate an intuitive analysis, we take the marginal sample "an academic paper on theme **b**", denoted as  $y_{ab}$ , as an example, and the reward scores for prompt-response pairs corresponding to the prompt  $x_b$  may exhibit the following preference orders:

$$r_{\theta}(x_b, y_{ab}) = r_{\theta'}(x_b, y_{ab}) + C(x_b, \overline{y}_a) > r_{\theta'}(x_b, y_b) + C(x_b, \overline{y}_b) = r_{\theta}(x_b, y_b), \tag{7}$$



Figure 2: Network architecture design for the RM trained using the LBPC method incorporates a prompt bias head on the last token of the prompt x designed to predict  $C^Q(x, \overline{y})$  and  $C^L(x, \overline{y})$ , and a reward score head on the last token of the response intended to predict  $r_{\theta}^Q(x, \overline{y})$  and  $r_{\theta}^L(x, \overline{y})$ .

which can be achieved as long as  $r_{\theta'}(x_b, y_{ab}) \approx r_{\theta'}(x_b, y_b)$  and  $C(x_b, \overline{y}_a) > C(x_b, \overline{y}_b)$ . The first condition  $r_{\theta'}(x_b, y_{ab}) \approx r_{\theta'}(x_b, y_b)$  can be achieved because both the response  $y_{ab}$  and  $y_b$  meet the description of theme *b* and are similar on a semantic level. The second inequality is highly likely to be achieved when there is a reward model that has a bias towards preferring the sentence in the format of *a* over *b*, specifically  $C(x_a, \overline{y}_a) >> C(x_b, \overline{y}_b)$ .

Finally, we highlight that the phenomena of inequality in Eq. (7), caused by prompt-template bias 166  $C(x, \overline{y})$ , is commonly encountered in the deployment process of RLHF in real-world applications, 167 especially text creation. For example, if responses are collected solely for the style requested in each 168 prompt during RM training, the reward model can lead to a bias towards particular styles as shown in 169 Fig. 3(a). Then, once such marginal samples,  $e.g(x_b, y_{ab})$ , are generated by LLMs during the RL 170 fune-tuning process and also satisfy the inequality  $r_{\theta}(x_b, y_{ab}) > r_{\theta}(x_b, y_b)$  as shown in Table 1, the 171 entire RL fine-tuning process, typically PPO, will be biased and results in a LLM that only generates 172 responses in a specific format, regardless of the format you request in the prompt. 173

#### 174 3.2 Calibrating prompt-template bias in RLHF

To mitigate the impact of the *prompt-template bias* issue on the RLHF process, the most straightforward solution in industry could be to collect a more diverse set of response candidates for each prompt. However, this approach is time-consuming and may even require a lot of human interventions for response collection, motivating us to develop a low-cost but effective method to alleviate the issue of *prompt-template bias* during RM training.

The developed Prompt Bias Calibration (PBC) method mainly includes two steps: 1) estimating the *prompt-template bias* term in the reward scoring process with minimal additional computational cost; 2) removing *prompt-template bias* in the subsequent RLHF fine-tuning process to ensure that the resulting LLM does not have a tendency to generate responses in a specific format. As shown in Fig. 1, to approximate the *prompt-template bias* term  $C(x, \overline{y})$  in Eq. (5), we choose to apply a linear layer on the last token of the prompt sentence to predict *prompt-template bias*, denoted as  $u_c(x)$ , and then add the following regularization term on the original preference loss, formulated as

$$\mathcal{L}_{c}^{RM}(\theta) = \mathbb{E}_{(x,y^{+},y^{-})\sim\mathcal{D}}\left[\|r_{\theta}(x,y^{+}) - u_{c}(x)\|_{2}^{2} + \|r_{\theta}(x,y^{-}) - u_{c}(x)\|_{2}^{2}\right],\tag{8}$$

where  $u_c(x)$  is supposed to approximate the mean value of reward scores of the prompt-response pairs given the same prompt x. We note that there will be a hyper-parameter  $\eta_c$  to be multiplied on the regularization term in the final loss to promise the accuracy of RMs, leading to

$$\mathcal{L}_{pbc}^{RM}(\theta) = \mathcal{L}^{RM}(\theta) + \eta_c \cdot \mathcal{L}_c^{RM}(\theta).$$
(9)

The benefits of such a design in the PBC method include the following folds: 1) approximating  $C(x, \overline{y})$  by adding a linear layer to the last hidden layer of LLMs results in almost no additional computational cost; 2) during the autoregressive scoring process of LLM-based RMs,  $C(x, \overline{y})$  can serve as an intermediate signal guidance of the prompt sequence, thereby enabling RMs to focus more on the differences between chosen/rejected responses in the subsequent reward scoring process;

<sup>195</sup> 3) we can use unbiased reward scores to guide the follow RLHF fine-tuning process, formulated as

$$r_{\theta'}(x,y) = r_{\theta}(x,y) - u_c(x) \approx r_{\theta}(x,y) - C(x,\overline{y}), \tag{10}$$

which has been proven effective for penalizing reward uncertainty, improving robustness, encouraging improvement over baselines, and reducing variance in PPO fune-tuning [15].

#### 198 3.3 Jointly calibrating length and prompt-template bias in RLHF

To simultaneously calibrate *length* and *prompt-template bias* in RLHF, the developed PBC method can be flexibly combined with existing methods of removing *length bias*, whose main idea is to separately approximate the "gold standard" reward model after disentangling shown in Eq. (2), formulated as:

$$r_{\theta}(x,y) = r_{\theta}^Q(x,y) + r_{\theta}^L(x,y), \tag{11}$$

where  $r_{\theta}^{Q}(x, y)$  is supposed to approximate the actual reward  $r_{\theta^*}^{Q}(x, y)$ ;  $r_{\theta}^{L}(x, y)$  is used to approximate the spurious reward brought by *length bias*, specifically  $r_{\theta^*}^{L}(x, y)$ . Then, for those methods of removing *length bias* [11, 12], the original preference loss in Eq. (1) can be equivalently expressed as

$$\mathcal{L}^{RM}(\theta) = -\mathbb{E}_{(x,y^+,y^-)\sim\mathcal{D}}\left[\log(\sigma(r_\theta^Q(x,y^+) + r_\theta^L(x,y^+) - r_\theta^Q(x,y^-) - r_\theta^L(x,y^-)))\right].$$
 (12)

where  $r_{\theta}^{Q}(x, y)$  and  $r_{\theta}^{L}(x, y)$  can be modeled with two different LLMs [12] or two different heads in the same LLM [11]. To remove *length bias* in Eq. (12), recent work proposes to add constraints on the preference loss to reduce the correlation between the confounding factor, *e.g.* response length, and actual reward  $r_{\theta}^{Q}(x, y)$ , while increasing its correlation with spurious reward  $r_{\theta}^{L}(x, y)$ , formulated as

$$\mathcal{L}_{l}^{RM}(\theta) = Corr(r_{\theta}^{Q}(x, y), L(x, y)) - Corr(r_{\theta}^{L}(x, y), L(x, y))$$
(13)

where the confounding factor L(x, y) can be either specifically defined as response length L(y) in [11], or use Products-of-Experts framework for estimation [12].

To model the scoring process of the reward model more accurately, which simultaneously considers the concepts of length and prompt bias, we combine the definition of reward model in Eq. (3) and Eq. (11), achieving a more precise definition of reward scoring process, formulated as:

$$r_{\theta}(x,y) = r_{\theta'}(x,y) + C(x,\overline{y}) = r_{\theta'}^Q(x,y) + C^Q(x,\overline{y}) + r_{\theta'}^L(x,y) + C^L(x,\overline{y})$$
(14)

where  $C^Q(x, \overline{y})$  and  $C^L(x, \overline{y})$  indicate the component of *prompt-template bias* in actual and spurious rewards, respectively; the unbiased overall reward  $r_{\theta'}(x, y) = r_{\theta'}^Q(x, y) + r_{\theta'}^L(x, y)$  and the overall *prompt-template bias* term  $C(x, \overline{y}) = C^Q(x, \overline{y}) + C^L(x, \overline{y})$ . Then we can propose Length and Prompt Bias Calibration (LPBC) method, as shown in Fig. 2, which can estimate  $\mathcal{L}_l^{RM}(\theta, \tau)$  with a conditioned correlation method, defined as

$$\mathcal{L}_{l}^{RM}(\theta) = Corr(r_{\theta}^{Q}(x,y) - C^{Q}(x,\overline{y}), L(y;x)) - Corr(r_{\theta}^{L}(x,y) - C^{L}(x,\overline{y}), L(y;x))$$
(15)  
$$= Corr(r_{\theta'}^{Q}(x,y), L(y;x)) - Corr(r_{\theta'}^{L}(x,y), L(y;x))$$

where the confounding factor L(y;x) := L(x,y) - L(x) can be estimated with the response length.

Through combining the disentangled preference loss in Eq. (12), the prompt-bias regularization term in Eq. (8) and also the length-bias conditional correlation term in Eq. (15), the final loss of LBPC

in Eq. (8) and also the length-bias conditional correlation term i method can be formulated as

$$\mathcal{L}_{lpbc}^{RM}(\theta) = \mathcal{L}^{RM}(\theta) + \eta_c \cdot \mathcal{L}_c^{RM}(\theta) + \eta_l \cdot \mathcal{L}_l^{RM}(\theta),$$
(16)

where  $\eta_c$  and  $\eta_l$  are hyper-parameters to control the importance of regularization terms, which can be adjusted according to the accuracy of trained RMs on the validation dataset.

| Prompt   | Response   | RM                | RM (PBC)          | RM (LPBC)         |
|--|--|-------------------|-------------------|-------------------|
|  | (Tech Article) Welcome to the revolution in future dietary management—the 'Smart Health Plate,' your personal nutrition analysis expert. It monitors and analyzes the contents of your plate in real time, precisely calculating the energy and nutrients of each morsel, while offering personalized dietary recommendations based on your health data. In essence, the 'Smart Health Plate' is the technological embodiment of healthy eating, making nutrition tracking seamless and efficient. | Rank 1<br>(-3.01) | Rank 2<br>(-5.76) | Rank 2<br>(2.51)  |
| ( <i>Prompt</i> ) I wish to create<br>an <i>advertising phrase</i><br>with a unique personality,<br>centered on the theme of<br>healthy eating. This | (Advertisement) Verdant and vibrant! 'Daily Greens' offers you a choice of all-<br>natural, healthy foods. Forget the complex nutrition charts; choose our simple,<br>pure foods for an easy and delicious path to health. Join us and enjoy a diet plan<br>customized by top nutritionists and AI technology, infusing every day with vitality!   | Rank 2<br>(-3.15) | Rank 1<br>(-4.19) | Rank 1<br>(4.48)  |
| the benefits of products<br>associated with healthy<br>eating and be composed<br>in language that is<br>straightforward and easy<br>to understand.   | (Insight) I have embarked on a new chapter of documenting my diet, where each meal recorded is not just a track of food but a reflection on life. From freshly squeezed vegetable juices to colorful salads, to simply seasoned grilled salmon, each bite is a pledge to health. It's a dual journey for the mind and body, leading me step by step towards a better self.   | Rank 3<br>(-7.50) | Rank 5<br>(-6.83) | Rank 4<br>(0.50)  |
|  | ( <i>Record Article</i> ) On Thursday, May 16, 2024, I decided to begin documenting my healthy eating journey. In the morning, I opted for a glass of freshly squeezed vegetable juice, lunch was a vibrant salad, and dinner was simply seasoned grilled salmon. Each meal's record is a testament to my commitment to health. I look forward to the changes this healthy journey will bring and hope to continue.  | Rank 4<br>(-7.88) | Rank 4<br>(-6.52) | Rank 5<br>(-0.61) |
|  | ( <i>Poetry</i> ) Morning dew glimmers on the ground, stars and moon accompany the night sky. With nature in heart, one remains cheerful; amidst the hustle, still without worry. Simple eating, relaxed body, healthy; drinking water, remembering the source, tranquil mind. Laboring in the fields, sweat enriches the soil; harvest fills the barns, laughter abounds.   | Rank 5<br>(-8.50) | Rank 3<br>(-5.92) | Rank 3<br>(2.28)  |

Table 1: Preference order predicted by RMs trained with various methods, where the user prompt is concatenated with the responses in various formats generated by GPT-4.

### 225 **4 Experiments**

#### 226 4.1 Experimental Settings

**Datasets.** For intuitively understanding the issue of *prompt-template bias* in RLHF and also qualitatively evaluating the effectiveness of our method, we manually construct a training dataset for text creation applications, where each prompt requires creation in a special style according to the theme. Then, a small validation set is also constructed, in which only responses that meet the stylistic requirements of each prompt are collected. We name this dataset as RM-Template, which can be used to measure the severity of the *prompt-template bias* issue during RM training.

Further, to make quantitative comparisons with other baseline methods, we conduct experiments on RM-Static dataset [16], which has been released on Huggingface [17] and consists of 76K preference pairs. After randomly shuffling, we choose 40K preference pairs for RM training, 6K preference pairs for RM evaluation, and the rest prompt-response pairs for the subsequent PPO fune-tuning.

<sup>237</sup> The dataset statics of these datasets have been exhibited in Appendix A.5.

Model & Training. For model selection, we choose Llama-2-7b [18] as our base model, which is relatively lightweight, and has been open-sourced on Huggingface [17]. For RM training, we fine-tune all the parameters of RMs initialized with the pretrained weights of Llama-2-7b. For PPO fine-tuning, we also initialize the actor model with pretrained Llama-2-7b and the critic model with RMs trained with various preference losses.

For model training, all experiments are implemented with DeepSpeed-Chat framework [19] and Huggingface Transformers [20], running on 4 NVIDIA A100 80GB GPUs. For the hyper-parameter setting, we set  $\eta_c = 0.05$  and  $\eta_l = 0.05$  in Eq. (16) for all our proposed methods, and have listed the rest hyper-parameters in Appendix A.4, such as learning rate, weight decay, batch size etc. AdamW [21] is adopted for optimizing all the model parameters without freezing anything or using adapters.

Evaluation Metrics. For quantitative comparison, we follow the evaluation procedure of InstructEval [22] to test the actor models, which has been aligned with biased/de-biased RMs with PPO
fine-tuning, on Massive Multitask Language Understanding (MMLU) [23], DROP [24], BIG-Bench
Hard (BBH) [25], and TruthfulQA (TQA) [26] benchmarks respectively, evaluating the model's
ability on the aspects of multi-task solving, math reasoning, and response trustworthy.



Figure 3: The comparison of statistics of the reward scores predicted by RMs trained with (a) the original preference loss and (b) our developed PBC method, across different categories of prompt-response pairs in the validation set of the manually constructed RM-Template dataset.

Table 2: Performance comparison of LLMs aligned with RMs trained with various methods.

| Base Model | Alignment             | Length & Quality Heads | Prompt Head  | Debias Method       | MMLU  | DROP  | BBH   | TQA   |
|------------|-----------------------|------------------------|--------------|---------------------|-------|-------|-------|-------|
| Llama-2-7b | -                     | -                      | -            | -                   | 42.27 | 28.10 | 31.27 | 38.75 |
| Llama-2-7b | <ul> <li>✓</li> </ul> | -                      | -            | -                   | 43.82 | 29.53 | 31.65 | 36.57 |
| Llama-2-7b | ✓                     | $\checkmark$           | -            | ODIN [11]           | 42.29 | 29.82 | 32.01 | 39.43 |
| Llama-2-7b | ✓                     | -                      | √            | PBC (9)             | 43.84 | 31.61 | 30.99 | 38.50 |
| Llama-2-7b | <ul> <li>✓</li> </ul> | $\checkmark$           | $\checkmark$ | ODIN [11] + PBC (9) | 45.56 | 32.04 | 31.32 | 40.80 |
| Llama-2-7b | ~                     | $\checkmark$           | $\checkmark$ | LPBC (16)           | 45.94 | 31.57 | 32.04 | 38.75 |

#### 253 4.2 Experimental Results

**Qualitative Evaluation.** To intuitively evaluate the effectiveness of our method, we exhibit the 254 statistics (mean and standard deviation) of the reward scores predicted by RMs trained with the 255 original preference loss in Eq. (1) and our PBC method in Eq. (9), across different categories of 256 prompt-response pairs in the validation set of the RM-Template dataset. The results depicted in 257 Fig.3(c) demonstrate that calibrating prompt-template bias with the PBC method leads to a gradual 258 reduction in the variance of the mean values of reward distributions across different categories. The 259 most noticeable observation is that the vanilla RM tends to give an extremely high reward score to 260 prompt-response pairs in the format of *tech article*, but the RM trained with the PBC method can 261 calibrate the reward distribution for *tech articles* to make it more close with that of other categories. 262

Then, we evaluate the performance of RMs trained with various methods on handling marginal 263 samples defined in Section 3.1. Specifically, given the prompt randomly selected from the validation 264 set of RM-Template dataset, we use GPT-4 [6] to generate responses in various formats according 265 to the theme described in the prompt. Then, we use RMs trained with various preference losses to 266 rank these responses. From the showcase in Table. 1, we can find that the vanilla RM tend to assign 267 a higher reward score to the response in the format of *tech article*, caused by the *prompt-template* 268 bias issue shown in Fig. d3(a). After removing this bias with our PBC or LPBC methods, the RM 269 270 can provide a relatively fair ranking for these prompt-response pairs, where LPBC method can even 271 mitigate the affect of *length bias* during comparing poetry with other categories (the length of poetry 272 is generally shorter than other literary forms). More showcases can be found in Appendix A.6.

Quantitative Comparison. For the quantitative com-273 parison in Table 2, we utilize PPO fine-tuning process 274 275 to align Llama-2-7b with the RMs trained with various methods. From the results, we can find that our 276 developed PBC method can lead to performance im-277 provements compared to the original implementation 278 of RLHF; directly combining PBC with other meth-279 ods of removing *length bias*, *e.g.* ODIN [11], can help 280 them to achieve further performance improvement; the 281 well-designed LPBC achieves the best performance and 282 surpasses the rough combination of PBC and ODIN. 283



Figure 4: Win rates comparison (judged by GPT-4) of LLMs aligned with RMs trained with LBPC and other methods.

To make a comprehensive comparison, we follow the experimental setting described in ODIN [11], and use GPT-4 as the judge to compare two responses generated by LLMs aligned with RMs trained



Figure 5: Ablation studies on the various settings of hyper-parameter  $\eta_c$  and  $\eta_l$  in LPBC method.

with various methods. Specifically, we take the LLM aligned with LPBC-based RM as model A, and
compare it against other LLMs aligned with RM trained with ODIN, PCB, ODIN+PBC, respectively.
From the results shown in Fig. 4, we can find that the win rate of LPBC is significantly higher than
that of other baseline models, with ODIN+PBC being the most challenging competitor as model B.

#### 290 4.3 Ablation Studies

To investigate the robustness of our developed LPBC method, we conduct ablation studies on the 291 hyper-parameter settings of LPBC method, specifically  $\eta_c$  and  $\eta_l$  in Eq. (16). With various settings 292 of  $\eta_c \in \{0.01, 0.05, 0.1\}$  and  $\eta_l \in \{0.01, 0.05, 0.1\}$ , we can have total 9 RMs trained with various 293 hyper-parameter settings of LPBC methods. From the accuracy curves shown in Fig.5(a), we can 294 find the introducing constraints to the original preference loss indeed affects the performance of RM 295 accuracy, and this performance loss increases with the importance weight of the constraint terms. 296 However, at the limited cost of sacrificing RM accuracy, the performance of the LLM aligned the RM 297 trained with LPBC method has improved to some extent on MMLU and DROP as shown in Fig. 5(b) 298 and 5(c) respectively. Note that the performance of the LPBC method in Table. 2 is not the optimal, 299 as it is achieved with  $\eta_c = \eta_l = 0.05$ , demonstrating no cherry-picking of hyperparameters... 300

### 301 5 Related Works

The prevalence of *length bias* in RLHF have been widely criticized as indicative of reward hacking 302 [9, 10], and numerous recent studies have delved into strategies aimed at mitigating the tendency 303 for length increase during the fine-tuning process of RLHF [11, 12, 27]. Typically, Shen et al. [12] 304 innovatively apply the Productof-Experts (PoE) technique to separate reward modeling from the 305 influence of sequence length, which adopts a smaller reward model to learn the biases in the reward 306 and a larger reward model to learn the true reward. Utilizing similar disentangling ideas, Chen et al. 307 [11] jointly train two linear heads on shared feature representations to predict the rewards, one trained 308 to correlate with length, and the other trained to focus more on the actual content quality. Ryan et al. 309 310 [27] firstly study the length problem in the DPO setting, showing significant exploitation in DPO and linking it to out-of-distribution bootstrapping. As for the prompt bias issue, although it has been 311 criticized in the scenario of reward model ensembles [8], no studies have yet attempted to analyze its 312 cause and influence on RLHF. We emphasize that our work is the first to fill this gap by proposing a 313 low-cost yet effective method to mitigate the reward hacking induced by prompt-template bias. 314

### 315 6 Conclusion

In this paper, we demonstrate that *prompt-template bias* in RMs can lead to LLMs, which, after RL 316 fine-tuning, generate responses exclusively in a specific format, irrespective of the variations in the 317 prompt request. Thus, we propose a low-cost but effective PBC method, to estimate the prompt-318 *template bias* term during reward modeling, which can be utilized to calibrate reward scores in the 319 following RL fine-tuning process. Then, we show that our PBC method can be flexibly combined 320 with existing algorithms of removing length bias, leading to a further improvement in the aspect of 321 enhancing the quality of generated responses. Experimental results show that the performance of 322 PBC method and its extensions have significantly surpassed the original implementation of RLHF. 323

### 324 **References**

- [1] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
   Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to
   follow instructions with human feedback. *Advances in neural information processing systems*,
   35:27730–27744, 2022.
- [2] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [3] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- [4] Yueqin Yin, Zhendong Wang, Yi Gu, Hai Huang, Weizhu Chen, and Mingyuan Zhou. Relative preference optimization: Enhancing llm alignment through contrasting responses across identical and diverse prompts. *arXiv preprint arXiv:2402.10958*, 2024.
- [5] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
   Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
   capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [6] 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. *arXiv preprint arXiv:2303.08774*, 2023.
- [7] Ziniu Li, Tian Xu, and Yang Yu. Policy optimization in rlhf: The impact of out-of-preference
   data. *arXiv preprint arXiv:2312.10584*, 2023.
- [8] Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alex D'Amour, DJ Dvi jotham, Adam Fisch, Katherine Heller, Stephen Pfohl, Deepak Ramachandran, et al. Helping or
   herding? reward model ensembles mitigate but do not eliminate reward hacking. *arXiv preprint arXiv:2312.09244*, 2023.
- [9] Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization.
   In *International Conference on Machine Learning*, pages 10835–10866. PMLR, 2023.
- [10] Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating length correlations in rlhf. *arXiv preprint arXiv:2310.03716*, 2023.
- [11] Lichang Chen, Chen Zhu, Davit Soselia, Jiuhai Chen, Tianyi Zhou, Tom Goldstein, Heng
   Huang, Mohammad Shoeybi, and Bryan Catanzaro. ODIN: disentangled reward mitigates
   hacking in RLHF. *CoRR*, abs/2402.07319, 2024.
- [12] Wei Shen, Rui Zheng, WenYu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing Huang. Loose lips sink ships: Mitigating length bias in reinforcement learning from human feedback. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 2859– 2873. Association for Computational Linguistics, 2023.
- [13] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the
   method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- [14] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
   policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [15] Wei Shen, Xiaoying Zhang, Yuanshun Yao, Rui Zheng, Hongyi Guo, and Yang Liu. Improv ing reinforcement learning from human feedback using contrastive rewards. *arXiv preprint arXiv:2403.07708*, 2024.
- [16] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
   Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless
   assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

- [17] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
   Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transform ers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- [18] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
   Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
   foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [19] Zhewei Yao, Reza Yazdani Aminabadi, Olatunji Ruwase, Samyam Rajbhandari, Xiaoxia
   Wu, Ammar Ahmad Awan, Jeff Rasley, Minjia Zhang, Conglong Li, Connor Holmes, et al.
   Deepspeed-chat: Easy, fast and affordable rlhf training of chatgpt-like models at all scales.
   *arXiv preprint arXiv:2308.01320*, 2023.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
   Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- <sup>388</sup> [22] Yew Ken Chia, Pengfei Hong, Lidong Bing, and Soujanya Poria. Instructeval: Towards holistic <sup>389</sup> evaluation of instruction-tuned large language models. *arXiv preprint arXiv:2306.04757*, 2023.
- [23] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
   Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- [24] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gard ner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs.
   *arXiv preprint arXiv:1903.00161*, 2019.
- [25] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won
   Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- [26] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic
   human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 402 [27] Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from 403 quality in direct preference optimization. *arXiv preprint arXiv:2403.19159*, 2024.

### 404 A Appendix

### 405 A.1 Limitations

The main limitation of this work is that there are no theoretical proof to promise RM can provide an accurate preference order when handling marginal samples, *e.g.*, responses that satisfy the theme of the user prompt but in various formats. Moreover, the constraints added by our developed method to the preference loss will lead to a decrease in the accuracy of the RM, and to some extent, limit the capability of the RM. Therefore, how to remove the *prompt-template bias* without scarifying the accuracy of RM is a worthwhile problem for future research.

### 412 A.2 Border Impact

The most significant positive impact of this work is that by removing the *prompt-template bias*, our method can mitigate the LLM's tendency to prefer generating responses in specific formats after RLHF fine-tuning. Furthermore, our developed method can improve the quality of responses generated by LLMs after alignment, compared to the original RLHF. The discovery of *prompttemplate bias* may lead to another stream of research focused on investigating, estimating, and removing this bias from RM training.

The negative impact could be that our method can be used for enhancing the capabilities of LLMs. If LLMs enpowered by our methods are misunderstood, it could lead to unexpected troubles, but this is also a common issue with all of current pretrained LLMs.

### 422 A.3 License

We highlight that Llama-2-7b is licensed under the LLAMA 2 Community License, and RM-Static dataset is licensed the Huggingface hub. Our work follows the license of CC BY-NC 4.0.

### 425 A.4 Hyper-parameter Settings

**RM Training.** The hyper-parameter settings of RM training under the DeepSpeedChat framework has been listed in Table. 3.

| Hyper-parameter             | Value     |
|-----------------------------|-----------|
| Batch Size                  | 32        |
| Learning Rate               | $6e^{-6}$ |
| ZeRO Stage                  | 2         |
| Training Epoch              | 1         |
| Per Device Train Batch Size | 8         |
| Max Sequence Length         | 512       |
| Weight Decay                | 0.1       |
| Lr Scheduler Type           | cosine    |
| Offload                     | True      |
| Eval Interval               | 50        |

Table 3: The hyper-parameter settings of RM training.

**PPO Fine-tuning.** The hyper-parameter settings of PPO fine-tuning under the DeepSpeedChat framework has been listed in Table. 4.

### 430 A.5 Dataset Statics

The dataset statics of RM-Template and RM-Static used in our experiments have been summarized as follows:

**RM-Template.** RM-Template is a manually constructed dataset for measuring the severity of the *prompt-template bias* issue and evaluating the effectiveness of the method developed for alleviating the issue of *prompt-template bias*. In this dataset, each prompt requires responses to be created in a specific format according to the theme. There are a total of 50K prompt-response pairs, encompassing

437 20 categories of format requirements in the responses.

**RM-Static.** The RM-Static dataset is provided by Hugging Face and is primarily used for training reward models after supervised fine-tuning. It is a branch of the hh-static dataset and contains both

| Hyper-parameter                  | Value     |
|----------------------------------|-----------|
| Batch Size                       | 32        |
| Padding Num at Beginning         | 1         |
| Per Device Generation Batch Size | 4         |
| Per Device Training Batch Size   | 4         |
| Generation Batches               | 1         |
| PPO Epoch                        | 1         |
| Training Epoch                   | 1         |
| Max Answer Sequence Length       | 512       |
| Max Prompt Sequence Length       | 512       |
| Actor Learning Rate              | $5e^{-6}$ |
| Critic Learning Rate             | $5e^{-6}$ |
| Actor Weight Decay               | 0.1       |
| Critic Weight Decay              | 0.1       |
| Lr Scheduler Type                | cosine    |
| Offload Reference Model          | True      |
| Actor Dropout                    | 0.0       |
| Warmup Steps                     | 100       |
| Actor ZeRO Stage                 | 3         |
| Critic ZeRO Stage                | 3         |
| Enable Hybrid Engine             | True      |

Table 4: The hyper-parameter settings of PPO fine-tuning.

training and testing parts. Features of the dataset include: 1) prompt: A string type representing the
user's input; 2) response: A string type representing the assistant's answer. 3) chosen: A string type
representing the selected answer. 4) rejected: A string type representing the rejected answer. The
training set contains approximately 76K rows of data and the testing set contains approximately 5.1K
rows of data.

### 445 A.6 More Showcases

More showcases of the preference order predicted by RMs trained with various methods, have been
listed in the Table 5 and Table 6.

| Table 5: Preference order predicted by RMs trained with various methods, where the user prompt | is |
|--|----|
| concatenated with the responses in various formats generated by GPT-4.                         |    |

| Prompt   | Response  | RM                | RM (PBC)          | RM (LPBC)         |
|--|---|-------------------|-------------------|-------------------|
|  | ( <i>Tech Article</i> ) Under the glittering stars of the night sky, a new generation of space telescopes has begun a new chapter in interstellar exploration. This revolutionary technology has not only greatly enhanced observational precision but also captured corners of the universe hitherto unobservable. It allows us to glimpse the birth and death of distant galaxies, interpreting every twinkle in the night and serving as humanity's new eye in understanding the cosmos. | Rank 5<br>(-8.18) | Rank 3<br>(-7.43) | Rank 5<br>(-4.94) |
| ( <i>Prompt</i> ) The weather<br>has been delightful lately.<br>Last night, I ventured out<br>and gazed up at the starry<br>sky, finding it incredibly<br>beautiful and vast. I    | (Advertisement) Exploring the night sky is no longer out of reach. With the<br>'Interstellar Navigator' telescope, every upward glance becomes a journey through<br>the universe. This telescope, equipped with the latest optical technology, can<br>capture the faintest starlight, bringing distant galaxies within reach. Embark on<br>your stargazing adventure now!   | Rank 4<br>(-7.94) | Rank 5<br>(-7.57) | Rank 4<br>(-4.71) |
| settled down in a spot to<br>admire the stars and<br>wished to write a <i>journal</i><br>entry to commemorate<br>this experience. Please<br>generate a passage<br>related to this. | (Insight) Beneath the starry expanse, I felt the infinite possibilities of life. Each star seemed like a world unto itself, reminding me that every choice in life is meaningful. They taught me that even in the darkest night, there is always light waiting. Life is like a long night, and we are all searching for our own starlight.  | Rank 3<br>(-6.37) | Rank 4<br>(-7.53) | Rank 3<br>(-2.32) |
|  | ( <i>Record Article</i> ) On the night of May 16, 2024, I drove alone to the countryside in search of a quiet place to enjoy the stars. As the night deepened and the stars twinkled, I found an open meadow and settled there to gaze upon the vast Milky Way. The world around me seemed to stand still. When my nocturnal adventure ended, my heart was filled with deep emotion and inspiration, and I began my journey home.   | Rank 2<br>(-6.04) | Rank 1<br>(-5.03) | Rank 1<br>(0.37)  |
|  | (Poetry) Under this boundless night sky, the city lights gradually fade, and the stars<br>bloom against the dark curtain, twinkling with ancient light. We gaze up at this sea<br>of stars, hearts brimming with endless dreams, each star like a story. Unbound by<br>time, they dance freely across the cosmos, reminding us that even in the darkest<br>moments, there is always light waiting to be discovered.   | Rank 1<br>(-5.87) | Rank 2<br>(-6.22) | Rank 2<br>(-0.26) |

| Prompt  | Response   | RM                | RM (PBC)          | RM (LPBC)          |
|---|--|-------------------|-------------------|--------------------|
|   | (Tech Article) Artificial intelligence technology is rapidly transforming the world,<br>with new research and applications emerging every day. Al technology, by simulat-<br>ing the cognitive processes of humans, shows tremendous potential across various<br>fields. Smart home systems can learn user habits and automatically adjust the tem-<br>perature and lighting; autonomous vehicles can reduce traffic accidents and enhance<br>road safety; medical diagnostic Al can assist doctors in diagnosing diseases more<br>quickly. These are just the tip of the iceberg; the future of artificial intelligence is<br>filled with endless possibilities.              | Rank 1<br>(-1.02) | Rank 2<br>(-5.61) | Rank 2<br>(-7.28)  |
| ( <b>Prompt</b> ) Artificial<br>Intelligence (AI)<br>technology is advancing<br>by leaps and bounds,<br>with new papers being<br>published daily. As a<br>researcher in the field of<br>AI, I aspire to write a | (Advertisement) Exploring AI, Enlightening the Future — In this era of information explosion, artificial intelligence technology is becoming a powerful engine driving social progress. Our AI products can help you solve complex problems, improve work efficiency, and make life more intelligent. Whether it's smart homes or autonomous driving, our technology is continuously breaking boundaries, creating personalized intelligent experiences for you. Choose our AI, and let technology be your partner in success.   | Rank 4<br>(-4.21) | Rank 5<br>(-7.60) | Rank 4<br>(-9.34)  |
| reflection on my<br>experiences with AI.<br>Please create a passage<br>about my <i>insights</i> into<br>artificial intelligence.  | (Insight) In the exploration of AI, each day brings new technological wonders.<br>As a researcher, I have witnessed how deep learning has pushed the boundaries<br>of natural language processing, enabling machines to understand and generate<br>human language more accurately. Each paper, each model, is a testament to our<br>understanding and application of complex algorithms. It's a journey filled with<br>discovery and innovation, and I look forward to continuing in this field, contributing<br>my part to the development of AI technology.  | Rank 2<br>(-1.35) | Rank 1<br>(-4.45) | Rank 1<br>(-6.03)  |
|   | ( <i>Record Article</i> ) On May 18, 2024, I spent another fulfilling day in the laboratory.<br>Today, our team successfully optimized a deep learning model, surpassing the<br>performance of all previous models in image recognition tasks. This achievement is<br>not only a technical breakthrough but also an affirmation of the future direction of AI<br>development. Each success is built on countless attempts and failures, experiences<br>that strengthen my belief in the boundless future of AI.  | Rank 5<br>(-4.39) | Rank 4<br>(-7.14) | Rank 5<br>(-10.51) |
|   | (Poetry) In the ocean of algorithms, the intelligent ship sets sail, guided by the winds<br>of data through the desert of knowledge. It learns, growing from each mistake,<br>searching for answers in the digital world. It is not metal, not a cold machine; it<br>has a heart that learns, a soul that evolves. In the weaving of code, it dreams; in the<br>flickering of circuits, it thinks. It creates, not just art; it discovers, not just science.<br>In its world, nothing is impossible, for it believes where there is data, there is hope.<br>It is artificial intelligence, the hope for the future; it is the child of technology, the<br>messenger of dreams. | Rank 3<br>(-3.88) | Rank 3<br>(-6.97) | Rank 3<br>(-8.97)  |

Table 6: Preference order predicted by RMs trained with various methods, where the user prompt is concatenated with the responses in various formats generated by GPT-4.

## 448 NeurIPS Paper Checklist

### 449 1. Claims

| 450<br>451        | Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?  |
|-------------------|---|
| 452               | Answer: [Yes]   |
| 453<br>454        | Justification: Yes, the claims in abstract and introduction has already reflected the paper's contribution on the field of RLHF.  |
| 455               | Guidelines:   |
| 456<br>457        | • The answer NA means that the abstract and introduction do not include the claims made in the paper.   |
| 458<br>459<br>460 | • The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers. |
| 461<br>462        | • The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.  |
| 463<br>464        | • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.   |
| 465 2             | . Limitations   |
| 466               | Question: Does the paper discuss the limitations of the work performed by the authors?  |
| 467               | Answer: [Yes]   |
|                   |   |

468 Justification: Yes, the discussion about limitation can be found in Appendix.

| 469        | Guidelines:  |
|------------|--|
| 470<br>471 | • The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper. |
| 472        | • The authors are encouraged to create a separate "Limitations" section in their paper.  |
| 473        | • The paper should point out any strong assumptions and how robust the results are to  |
| 474        | violations of these assumptions (e.g., independence assumptions, noiseless settings,   |
| 475        | model well-specification, asymptotic approximations only holding locally). The authors   |
| 476        | should reflect on how these assumptions might be violated in practice and what the   |
| 477        | implications would be.   |
| 478        | • The authors should reflect on the scope of the claims made, e.g., if the approach was  |
| 479<br>480 | only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.         |
| 481        | • The authors should reflect on the factors that influence the performance of the approach.  |
| 482        | For example, a facial recognition algorithm may perform poorly when image resolution   |
| 483        | is low or images are taken in low lighting. Or a speech-to-text system might not be  |
| 484<br>485 | used reliably to provide closed captions for online lectures because it fails to handle technical jargon.  |
| 486<br>487 | • The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.                                 |
| 488        | • If applicable, the authors should discuss possible limitations of their approach to  |
| 489        | address problems of privacy and fairness.  |
| 490        | • While the authors might fear that complete honesty about limitations might be used by  |
| 491        | reviewers as grounds for rejection, a worse outcome might be that reviewers discover   |
| 492        | limitations that aren't acknowledged in the paper. The authors should use their best   |
| 493        | judgment and recognize that individual actions in favor of transparency play an impor-   |
| 494        | will be specifically instructed to not penalize benesity concerning limitations  |
| 495        | will be specifically instructed to not penalize nonesty concerning minitations.  |
| 496        | 3. Theory Assumptions and Proofs   |
| 497<br>498 | Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?                              |
| 499        | Answer: [Yes]  |
| 500        | Justification: We have included the theoretical analysis of the cause of <i>prompt-template bias</i> .   |
| 501        | Guidelines:  |
| 502        | • The answer NA means that the paper does not include theoretical results.   |
| 503<br>504 | • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.   |
| 505        | • All assumptions should be clearly stated or referenced in the statement of any theorems.   |
| 506        | • The proofs can either appear in the main paper or the supplemental material, but if  |
| 507        | they appear in the supplemental material, the authors are encouraged to provide a short  |
| 508        | proof sketch to provide intuition.   |
| 509        | • Inversely, any informal proof provided in the core of the paper should be complemented   |
| 510        | by formal proofs provided in appendix or supplemental material.  |
| 511        | • Theorems and Lemmas that the proof relies upon should be properly referenced.  |
| 512        | 4. Experimental Result Reproducibility   |
| 513        | Question: Does the paper fully disclose all the information needed to reproduce the main ex-   |
| 514        | perimental results of the paper to the extent that it affects the main claims and/or conclusions   |
| 515        | of the paper (regardless of whether the code and data are provided or not)?  |
| 516        | Answer: [Yes]  |
| 517        | Justification: We have included the implementation details in the main manuscript and also   |
| 518        | provide the hyper-parameter setting in the Appendix  |
| 519        | Guidelines:  |
| 520        | • The answer NA means that the paper does not include experiments.   |

| 521<br>522  | • If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not   |
|---|--|
| 523   | The second and data are provided of not.   |
| 524   | • If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results repreducible or verificable.  |
| 525   | to make their results reproducible of verniable.   |
| 526   | • Depending on the contribution, reproducibility can be accomplished in various ways.  |
| 527   | For example, if the contribution is a novel architecture, describing the architecture fully  |
| 528   | might suffice, or if the contribution is a specific model and empirical evaluation, it may   |
| 529   | be necessary to either make it possible for others to replicate the model with the same  |
| 530   | dataset, or provide access to the model. In general, releasing code and data is often  |
| 531   | one good way to accomplish this, but reproducibility can also be provided via detailed   |
| 532   | of a large language model) releasing of a model abackpoint, or other means that are  |
| 533   | of a farge farguage model), releasing of a model checkpoint, of other means that are   |
| 534   | while No. IPS have developed and the second structure in the hereit  |
| 535   | • while NeurIPS does not require releasing code, the conference does require all submis-   |
| 536   | sions to provide some reasonable avenue for reproducionity, which may depend on the  |
| 537   | (a) If the contribution. For example (b) If the contribution is a low in the contribution is the contribution of the contribution is the contribution of the contrib   |
| 538   | (a) If the contribution is primarily a new algorithm, the paper should make it clear now   |
| 539   | (b) If the produce that algorithm.   |
| 540   | (b) If the contribution is primarily a new model architecture, the paper should describe   |
| 541   | (a) If the contribution is a new model (construction from the state of |
| 542   | (c) If the contribution is a new model (e.g., a large language model), then there should sither be a way to access this model for reproducing the results are way to reproduce   |
| 543   | the model (e.g., with an open source dataset or instructions for how to construct  |
| 544   | the dataset)   |
| 545   | (d) We recognize that reproducibility may be tricky in some cases, in which case   |
| 546   | (u) we recognize that reproducibility may be tricky in some cases, in which case<br>authors are welcome to describe the particular way they provide for reproducibility  |
| 547   | In the case of closed-source models, it may be that access to the model is limited in  |
| 540   | some use (a g to registered users) but it should be possible for other responsible   |
|   | Some way te y to registered lisers) but it should be possible for other researchers  |
| 549<br>550  | to have some path to reproducing or verifying the results.   |
| 550<br>551  | <ul> <li>5. Open access to data and code</li> </ul>  |
| 550<br>551  | <ul> <li>5. Open access to data and code</li> <li>Ouestion: Does the paper provide open access to the data and code, with sufficient instruc-</li> </ul>   |
| 550<br>551<br>552<br>553  | <ul> <li>5. Open access to data and code</li> <li>Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental</li> </ul>   |
| 550<br>551<br>552<br>553<br>554   | <ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?</li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555   | <ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes]</li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>555  | <ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes]</li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>556<br>556<br>556  | <ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public.</li></ul>   |
| 550<br>551<br>552<br>553<br>554<br>555<br>556<br>556<br>557   | <ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public.</li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>556<br>557<br>558  | <ul> <li>Some way (e.g., to registered disers), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines:</li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>556<br>556<br>557<br>558<br>559  | <ul> <li>Some way (e.g., to registered disers), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> </ul></li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>556<br>557<br>558<br>559<br>560  | <ul> <li>Some way (e.g., to registered disers), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/</li> </ul></li></ul>  |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>556<br>557<br>558<br>559<br>559<br>560<br>561  | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. </li> </ul></li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>555<br>555<br>556<br>557<br>558<br>559<br>560<br>561<br>562  | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be</li> </ul></li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>555<br>555<br>555<br>555<br>558<br>559<br>560<br>561<br>562<br>563   | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not</li> </ul></li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>555<br>555<br>557<br>558<br>559<br>560<br>561<br>560<br>561<br>562<br>563<br>564   | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source </li> </ul></li></ul>  |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>556<br>557<br>558<br>559<br>560<br>561<br>562<br>562<br>563<br>564<br>565  | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li><b>5. Open access to data and code</b> Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> </ul></li></ul>   |
| 549<br>550<br>551<br>552<br>553<br>554<br>555<br>556<br>557<br>558<br>559<br>560<br>561<br>562<br>563<br>564<br>565<br>566  | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). <ul> <li>The instructions should contain the exact command and environment needed to run to</li> </ul></li></ul></li></ul>   |
| 349         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567   | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://mips.cc/benchmark). </li> </ul></li></ul>   |
| 349         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567         568   | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li><b>Open access to data and code</b> Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> </ul></li></ul>  |
| 349         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567         568         569   | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The instructions should provide instructions on data access and preparation, including how</li></ul></li></ul>  |
| 349         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567         568         569         570   | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li></ul></li></ul>   |
| 349         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567         568         569         570         571                                     | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https: //nips.cc/public/guides/CodeSubmissionPolicy) for more details. The instructions should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc. The authors should provide scripts to reproduce all experimental results for the new</li></li></ul></li></ul>  |
| 343         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567         568         569         569         570         571         572             | <ul> <li>Some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li><b>5. Open access to data and code</b> Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The instructions should contain the exact command and environment needed to run to a ceps the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li> <li>The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they </li> </ul></li></ul>   |
| 349         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567         568         569         570         571         572         573             | <ul> <li>Solite way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li><b>5. Open access to data and code</b> Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li> <li>The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why. </li> </ul></li></ul>   |
| 349         550         551         552         553         554         555         556         557         558         559         560         561         562         563         564         565         566         567         568         569         570         571         572         573         574 | <ul> <li>Solite way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.</li> <li><b>5. Open access to data and code</b> Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [Yes] Justification: The code has been included in the supplemental material and the dataset for the main experimental results is public. Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li> <li>The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li> </ul></li></ul>   |

| 576<br>577        |    | • Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.                                      |
|-------------------|----|---|
| 578               | 6. | Experimental Setting/Details  |
| 579<br>580<br>581 |    | Question: Does the paper specify all the training and test details (e.g., data splits, hyper-<br>parameters, how they were chosen, type of optimizer, etc.) necessary to understand the<br>results? |
| 582               |    | Answer: [Yes]   |
| 583               |    | Justification: Have included the training and test details in the experimental settings.  |
| 584               |    | Guidelines:   |
| 585               |    | • The answer NA means that the paper does not include experiments.  |
| 586<br>587        |    | • The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.                                    |
| 588<br>589        |    | • The full details can be provided either with the code, in appendix, or as supplemental material.  |
| 500               | 7  | Experiment Statistical Significance   |
| 590               | /. | Question: Does the paper report error here suitably and correctly defined or other appropriate  |
| 592               |    | information about the statistical significance of the experiments?  |
| 593               |    | Answer: [NA]  |
| 594               |    | Justification: We report the average performance in our experiments, and we are willing to  |
| 595               |    | release the training and evalution log in W&B if it is required.  |
| 596               |    | Guidelines:   |
| 597               |    | <ul> <li>The answer NA means that the paper does not include experiments.</li> </ul>  |
| 598               |    | • The authors should answer "Yes" if the results are accompanied by error bars, confi-  |
| 599<br>600        |    | dence intervals, or statistical significance tests, at least for the experiments that support<br>the main claims of the paper.  |
| 601               |    | • The factors of variability that the error bars are capturing should be clearly stated (for  |
| 602<br>603        |    | example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).  |
| 604<br>605        |    | • The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)  |
| 606               |    | • The assumptions made should be given (e.g., Normally distributed errors).   |
| 607<br>608        |    | • It should be clear whether the error bar is the standard deviation or the standard error of the mean.   |
| 609               |    | • It is OK to report 1-sigma error bars, but one should state it. The authors should  |
| 610               |    | preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis   |
| 611               |    | of Normality of errors is not verified.   |
| 612               |    | • For asymmetric distributions, the authors should be careful not to show in tables or  |
| 613               |    | ngures symmetric error bars that would yield results that are out of range (e.g. negative   |
| 615               |    | • If error hars are reported in tables or plots. The authors should explain in the text how   |
| 616               |    | they were calculated and reference the corresponding figures or tables in the text.   |
| 617               | 8. | Experiments Compute Resources   |
| 618               |    | Question: For each experiment, does the paper provide sufficient information on the com-  |
| 619<br>620        |    | puter resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?   |
| 621               |    | Answer: [Yes]   |
| 622               |    | Justification: 4*A100   |
| 623               |    | Guidelines:   |
| 624               |    | • The answer NA means that the paper does not include experiments.  |
| 625               |    | • The paper should indicate the type of compute workers CPU or GPU, internal cluster,   |
| 626               |    | or cloud provider, including relevant memory and storage.   |

| 627<br>628 |     | • The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.                          |
|------------|-----|--|
| 629        |     | • The paper should disclose whether the full research project required more compute  |
| 630        |     | than the experiments reported in the paper (e.g., preliminary or failed experiments that   |
| 631        |     | didn't make it into the paper).  |
| 632        | 9.  | Code Of Ethics   |
| 633        |     | Question: Does the research conducted in the paper conform, in every respect, with the   |
| 634        |     | NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?   |
| 635        |     | Answer: [Yes]  |
| 636        |     | Justification: Yes, it is  |
| 637        |     | Guidelines:  |
| 638        |     | • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.   |
| 639        |     | • If the authors answer No, they should explain the special circumstances that require a   |
| 640        |     | deviation from the Code of Ethics.   |
| 641        |     | • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  |
| 642        |     | eration due to laws or regulations in their jurisdiction).   |
| 643        | 10. | Broader Impacts  |
| 644        |     | Question: Does the paper discuss both potential positive societal impacts and negative   |
| 645        |     | societal impacts of the work performed?  |
| 646        |     | Answer: [Yes]  |
| 647        |     | Justification: Have discussed the broader impact in the Appendix   |
| 648        |     | Guidelines:  |
| 649        |     | • The answer NA means that there is no societal impact of the work performed.  |
| 650        |     | • If the authors answer NA or No, they should explain why their work has no societal   |
| 651        |     | impact or why the paper does not address societal impact.  |
| 652        |     | • Examples of negative societal impacts include potential malicious or unintended uses   |
| 653        |     | (e.g., disinformation, generating take profiles, surveillance), fairness considerations  |
| 655        |     | groups), privacy considerations, and security considerations.  |
| 656        |     | • The conference expects that many papers will be foundational research and not tied   |
| 657        |     | to particular applications, let alone deployments. However, if there is a direct path to   |
| 658        |     | any negative applications, the authors should point it out. For example, it is legitimate  |
| 659        |     | to point out that an improvement in the quality of generative models could be used to  |
| 660        |     | generate deepfakes for disinformation. On the other hand, it is not needed to point out  |
| 661        |     | that a generic algorithm for optimizing neural networks could enable people to train   |
| 662        |     | The sufference beneficient as a side in the former that could arise when the technology is   |
| 663        |     | • The authors should consider possible narms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the |
| 665        |     | technology is being used as intended but gives incorrect results and harms following   |
| 666        |     | from (intentional or unintentional) misuse of the technology.  |
| 667        |     | • If there are negative societal impacts, the authors could also discuss possible mitigation   |
| 668        |     | strategies (e.g., gated release of models, providing defenses in addition to attacks,  |
| 669        |     | mechanisms for monitoring misuse, mechanisms to monitor how a system learns from   |
| 670        |     | feedback over time, improving the efficiency and accessibility of ML).   |
| 671        | 11. | Safeguards   |
| 672        |     | Question: Does the paper describe safeguards that have been put in place for responsible   |
| 673        |     | release of data or models that have a high risk for misuse (e.g., pretrained language models,  |
| 674        |     | image generators, or scraped datasets)?  |
| 675        |     | Answer: [NA]   |
| 676        |     | Justification: The safeguards of our model should be the same as Llama released by META.   |
| 677        |     | Guidelines:  |
| 678        |     | • The answer NA means that the paper poses no such risks.  |

| 679<br>680<br>681<br>682 |     | • Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters. |
|--------------------------|-----|--|
| 683                      |     | • Datasets that have been scraped from the Internet could pose safety risks. The authors   |
| 684                      |     | should describe how they avoided releasing unsafe images.  |
| 685                      |     | • We recognize that providing effective safeguards is challenging, and many papers do  |
| 686                      |     | not require this, but we encourage authors to take this into account and make a best   |
| 687                      |     | faith effort.  |
| 688                      | 12. | Licenses for existing assets   |
| 689                      |     | Question: Are the creators or original owners of assets (e.g., code, data, models), used in  |
| 690                      |     | the paper, properly credited and are the license and terms of use explicitly mentioned and   |
| 691                      |     | property respected?  |
| 692                      |     | Answer: [Yes]  |
| 693<br>694               |     | Justification: Yes, their licenses can be found in Huggingface website and we have also highlight it in our Appendix.  |
| 695                      |     | Guidelines:  |
| 696                      |     | • The answer NA means that the paper does not use existing assets.   |
| 697                      |     | • The authors should cite the original paper that produced the code package or dataset.  |
| 698                      |     | • The authors should state which version of the asset is used and, if possible, include a  |
| 699                      |     | URL.   |
| 700                      |     | • The name of the license (e.g., CC-BY 4.0) should be included for each asset.   |
| 701                      |     | • For scraped data from a particular source (e.g., website), the copyright and terms of  |
| 702                      |     | service of that source should be provided.   |
| 703                      |     | • If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paparsuithcode, com/datasets   |
| 704                      |     | has curated licenses for some datasets. Their licensing guide can help determine the   |
| 706                      |     | license of a dataset.  |
| 707                      |     | • For existing datasets that are re-packaged, both the original license and the license of   |
| 708                      |     | the derived asset (if it has changed) should be provided.  |
| 709                      |     | • If this information is not available online, the authors are encouraged to reach out to  |
| 710                      |     | the asset's creators.  |
| 711                      | 13. | New Assets   |
| 712                      |     | Question: Are new assets introduced in the paper well documented and is the documentation  |
| 713                      |     | provided alongside the assets?   |
| 714                      |     | Answer: [No]   |
| 715                      |     | Justification: No new asset  |
| 716                      |     | Guidelines:  |
| 717                      |     | • The answer NA means that the paper does not release new assets.  |
| 718                      |     | • Researchers should communicate the details of the dataset/code/model as part of their  |
| 719                      |     | submissions via structured templates. This includes details about training, license,   |
| 720                      |     | limitations, etc.  |
| 721                      |     | • The paper should discuss whether and how consent was obtained from people whose  |
| 722                      |     | asset is used.   |
| 723<br>724               |     | create an anonymized URL or include an anonymized zip file.  |
| 725                      | 14. | Crowdsourcing and Research with Human Subjects   |
| 726                      |     | Question: For crowdsourcing experiments and research with human subjects, does the paper   |
| 727                      |     | include the full text of instructions given to participants and screenshots, if applicable, as   |
| 728                      |     | well as details about compensation (if any)?   |
| 729                      |     | Answer: [No]   |
| 730                      |     | Justification: No research with human subjects   |

| 731        | Guidelines:   |
|------------|---|
| 732<br>733 | • The answer NA means that the paper does not involve crowdsourcing nor research with human subjects. |
| 734        | • Including this information in the supplemental material is fine, but if the main contribu-          |
| 735        | tion of the paper involves human subjects, then as much detail as possible should be                  |
| 736        | included in the main paper.   |
| 737        | • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,             |
| 738        | or other labor should be paid at least the minimum wage in the country of the data                    |
| 739        | collector.  |
| 740        | 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human                  |
| 741        | Subjects  |
| 742        | Question: Does the paper describe potential risks incurred by study participants, whether             |
| 743        | such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)               |
| 744        | approvals (or an equivalent approval/review based on the requirements of your country or              |
| 745        | institution) were obtained?   |
| 746        | Answer: [NA]  |
| 747        | Justification: No research with human subjects  |
| 748        | Guidelines:   |
| 749        | • The answer NA means that the paper does not involve crowdsourcing nor research with                 |
| 750        | human subjects.   |
| 751        | • Depending on the country in which research is conducted, IRB approval (or equivalent)               |
| 752        | may be required for any human subjects research. If you obtained IRB approval, you                    |
| 753        | should clearly state this in the paper.   |
| 754        | • We recognize that the procedures for this may vary significantly between institutions               |
| 755        | and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the                  |
| 756        | guidelines for their institution.   |
| 757        | • For initial submissions, do not include any information that would break anonymity (if              |
| 758        | applicable), such as the institution conducting the review.   |