

LIRE: listwise reward enhancement for preference alignment

Mingye Zhu¹, Yi Liu^{2*}, Lei Zhang¹
Junbo Guo² and Zhendong Mao¹

¹University of Science and Technology of China, Hefei, China

²State Key Laboratory of Communication Content Cognition, Beijing, China
mingyechu@mail.ustc.edu.cn, gavin1332@gmail.com
{leizh23, zdmao}@ustc.edu.cn, guojunbo@people.cn

Abstract

Recently, tremendous strides have been made to align the generation of Large Language Models (LLMs) with human values to mitigate toxic or unhelpful content. Leveraging Reinforcement Learning from Human Feedback (RLHF) proves effective and is widely adopted by researchers. However, implementing RLHF is complex, and its sensitivity to hyperparameters renders achieving stable performance and scalability challenging. Furthermore, prevailing approaches to preference alignment primarily concentrate on pairwise comparisons, with limited exploration into multi-response scenarios, thereby overlooking the potential richness within the candidate pool. For the above reasons, we propose a new approach: *Listwise Reward Enhancement for Preference Alignment* (LIRE), a gradient-based reward optimization approach that incorporates the offline rewards of multiple responses into a streamlined listwise framework, thus eliminating the need for online sampling during training. LIRE is straightforward to implement, requiring minimal parameter tuning, and seamlessly aligns with the pairwise paradigm while naturally extending to multi-response scenarios. Moreover, we introduce a self-enhancement algorithm aimed at iteratively refining the reward during training. Our experiments demonstrate that LIRE consistently outperforms existing methods across several benchmarks on dialogue and summarization tasks, with good transferability to out-of-distribution data, assessed using proxy reward models and human annotators.

1 Introduction

While a growing plethora of large language models (LLMs) have exhibited incredible performance in a broadening scope of tasks and applications such as summarization, machine translation, and dialog generation (Nakano et al., 2021; Stiennon et al., 2020; Brown et al., 2020; Zhao et al., 2023a), they

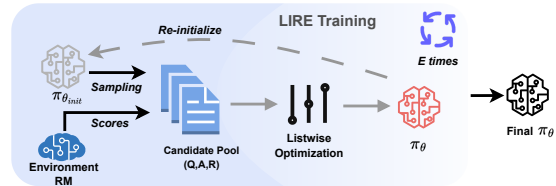


Figure 1: **Training pipeline of the proposed LIRE framework.** The candidate pool is initially constructed by gathering responses A with arbitrary policy $\pi_{\theta_{init}}$. Subsequently, the scored responses with their query are optimized in a listwise manner. The "dashed" line means it is optional to re-initialize the updated model π_{θ} as the sampling policy and generates fresh responses that substitute the prior ones within the candidate pool.

can still output contents that are harmful, biased or simply do not agree with standard human perception (Mathur et al., 2020; Fernandes et al., 2023). This is an inherent problem existing in the extensive data sources during model training (Ouyang et al., 2022; Bai et al., 2022b; Song et al., 2023), and can be alleviated by incorporating certain restrictions or limitations to align the output generation towards human desires and specifications (Ngo, 2022; Kenton et al., 2021).

Existing methods focus on employing Reinforcement Learning from Human Feedback (RLHF) to fine-tune the pre-trained LLMs (Christiano et al., 2017; Stiennon et al., 2020; Ouyang et al., 2022; Xue et al., 2023), which introduces a paradigm that involves leveraging supervised fine-tuning (SFT) on the initial models, fitting the reward model to human preferences, and then using Reinforcement Learning (RL) algorithms such as Proximal Policy Optimization (PPO) (Schulman et al., 2017) to optimize a policy that doesn't drift overly far from the original model.

However, PPO is optimized in a pointwise manner based on the sparse rewards, penalizing fragments within a sentence equally and disregarding the truly informative parts. Additionally, PPO requires online sampling during training, which im-

*Corresponding author: Yi Liu

pairs computational efficiency and stability. Alternatively, methods such as DPO (Rafailov et al., 2023) and RRHF (Yuan et al., 2023) primarily leverage a pairwise comparison between a positive and a negative sample and transform the RL problem into an offline setting. However, the performance of pairwise comparison is heavily dependent on the quality of the sample pairs, and trivial negatives may yield suboptimal results. Moreover, given a large candidate pool, performing pairwise comparisons among multiple samples entails a significant computational complexity.

To fully exploit the diversity among multiple samples and better identify high-reward segments within sentences, we propose an RL-free listwise approach: *Listwise Reward Enhancement for Preference Alignment* (LIRE). Unlike existing algorithms (Rafailov et al., 2023; Song et al., 2023) that employ the Bradley-Terry model (Bradley and Terry, 1952) or Plackett-Luce (PL) models (Plackett, 1975) to model the preferences, we reformulate the response probability distribution and implicitly model the preferences with the offline rewards. The probabilities of responses are updated under the guidance of the rewards, entailing a more accurate indication of how the responses are preferred than their binary format counterpart of "chosen" and "reject" labels.¹

Notably, the concept of *listwise* in LIRE is different from the traditional *listwise* approach proposed in the Learning-to-Rank literature (Cao et al., 2007; Xia et al., 2008), which is largely based on the PL model and requires a ranking permutation determined by the *position* given any labels. The key idea for the PL model is that the choice in the i -th position in a rank π only depends on the candidates not chosen at previous positions, leading to a time complexity of at least $\mathcal{O}(n)$ depending on different surrogate losses. In contrast, in LIRE we compute the exponential probability distribution only once, making it efficient to compute even if n increases.

The overarching concept is illustrated in Figure 1: we first gather responses A for queries Q from different sources, after which the rewards are collected before the optimization begins. An optional strategy to further boost performance is

¹In this sense, LIRE can also be viewed as an optimization objective under the paradigm of Minimum Bayes Risk (MBR) decoding (Bertsch et al., 2023), where the method for choosing the output is based not on the output with the highest probability, but the output with the lowest risk (highest reward) among multiple candidates.

re-initializing the trained model to generate fresh responses and conduct iterative training.

LIRE is easy to implement, free from heavy parameter tuning, and perfectly fitted within the pairwise paradigm with natural extension towards the multi-response scenario. Experiments of the state-of-the-art methods are fairly conducted on multiple benchmarks of dialogue and summarization tasks. The results show that the proposed LIRE achieves superior and consistent performance in all the experiments, using strong 7B-LLMs as base models and various proxy reward models, GPT-4 as well as human evaluators for assessment.

2 Preliminaries

In this section, we illustrate the motivation for the LIRE framework and the related preliminaries. We start with the objective function of the Policy Gradient (PG) methods:

$$J(\theta) = \sum_{\tau} P(\tau, \theta) R(\tau), \quad (1)$$

where $P(\tau, \theta)$ is the probability of the trajectory, and $R(\tau)$ is the corresponding return. This function can be formulated into a differentiable function using the PG Theorem:

$$\nabla J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla \log \pi_{\theta}(a_t | s_t) R(\tau)], \quad (2)$$

where a_t and s_t are the action and state at time step t of a certain trajectory τ . The ultimate goal of PG methods is to maximize the rewards of the trajectories under the policy π_{θ} .

Supposing we have a set of queries x and answers y , the optimization objective widely used in RLHF settings (Ouyang et al., 2022; Stiennon et al., 2020; Ziegler et al., 2019) is:

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} \left(r_{\phi}(x, y) \right) - \beta \mathbb{D}_{\text{KL}} \left(\pi_{\theta}(y|x) || \pi_{\text{ref}}(y|x) \right), \quad (3)$$

where r_{ϕ} is the well-trained reward function. The sampling efficiency is limited in this process since the training data has to be sampled online as policy π_{θ} updates, with KL penalty included to avoid utter policy drifting (Schulman et al., 2017).

To better approximate $P(\tau, \theta)$ and thus a better indicator of the expected rewards in Equation 1, we propose to leverage multiple pre-allocated responses to reformulate the trajectory probability

distribution. We also transform the RL learning problem into an offline setting, improving computational efficiency and stability.

3 Methodology

In this section, we first introduce the construction of the LIRE objective and its relation with another popular objective from a theoretical aspect. Subsequently, we introduce a self-enhancement algorithm that further boosts the model preference. Please note that directly training with pairwise preference data in a single stage (without Algorithm 1) will suffice to bring out impressive performance and only experiments in Section 4.7 leverages iterative sampling and training.

3.1 The construction of the LIRE objective

In this section, we reformulate the preference alignment problem and introduce our LIRE framework. Firstly, we assume a set of queries $\mathbf{Q} = \{\mathbf{x}^{(i)}\}$ is given, $i \in \{1, \dots, N\}$ and each query is associated with M responses $\mathbf{A}^{(i)} = \{\mathbf{y}_1^{(i)}, \dots, \mathbf{y}_M^{(i)}\}$. Furthermore, each response $\mathbf{y}_j^{(i)}$ for query $\mathbf{x}^{(i)}$ is paired with a score $R(\mathbf{x}^{(i)}, \mathbf{y}_j^{(i)})$ by some Reward Model RM.

Next, we omit the superscript of (i) for clarity. During training, we aim to learn a language model parameterized by θ , which generates responses of better alignment with human preferences. First, we define a set of token prediction probabilities conditioned on \mathbf{x} as $\mathbf{P}_{\pi_\theta}(\mathbf{y}_{j,k}|\mathbf{x}) \in \mathbb{R}^{L \times V}$, where L is the sequence length and V the vocabulary size. The probability of the sentence \mathbf{y}_j with K tokens takes the form in an autoregressive nature:

$$\pi_\theta(\mathbf{y}_j|\mathbf{x}) = \prod_{k=1}^K \mathbf{P}_{\pi_\theta}(\mathbf{y}_{j,k}|\mathbf{x}, \mathbf{y}_{j,<k}). \quad (4)$$

To exploit the diversity present in different responses, we reformulate the response probability distribution against the entire response set \mathbf{A} as:

$$P_{\pi_\theta}(\mathbf{y}|\mathbf{x}, \mathbf{A}) = \frac{\exp(\frac{1}{T} \log \pi_\theta(\mathbf{y}|\mathbf{x}))}{\sum_{j=1}^M \exp(\frac{1}{T} \log \pi_\theta(\mathbf{y}_j|\mathbf{x}))}, \quad (5)$$

where T is a temperature parameter to control the smoothness of the probability distribution. Equation 5 can be regarded as an updated approximation of the underlying response (trajectory) distribution,

and we next derive the listwise loss as:

$$\begin{aligned} J(\theta) &= -\mathbb{E}_{\mathbf{x} \sim q(\cdot)} \mathbb{E}_{\mathbf{y} \sim P_{\pi_\theta}(\cdot|\mathbf{x}, \mathbf{A})} R(\mathbf{x}, \mathbf{y}) \\ &= -\mathbb{E}_{\mathbf{x} \sim q(\cdot)} \sum_{j=1}^M P_{\pi_\theta}(\mathbf{y}_j|\mathbf{x}, \mathbf{A}) R(\mathbf{x}, \mathbf{y}_j), \end{aligned} \quad (6)$$

where $q(\cdot)$ denotes the distribution of the queries. In practice, we apply softmax to the reward scores of a single query due to its property of translation invariance. By doing so we mitigate the influence of different reward scales and maintain stable training parameter settings. Next, to develop a general perception of what the model learns through the process, we illustrate the gradient of $J(\theta)$, whose derivation process can be found in Appendix A.1:

$$\begin{aligned} \nabla_\theta J(\theta) &= -\frac{1}{T} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y} \sim \pi_\theta(\cdot|\mathbf{x})} \left[\frac{\nabla \pi_\theta(\mathbf{y}|\mathbf{x})}{\pi_\theta(\mathbf{y}|\mathbf{x})} \right. \\ &\quad \left. \times \left(R(\mathbf{x}, \mathbf{y}) - \mathbb{E}_{(\mathbf{y}' \sim \pi_\theta(\cdot|\mathbf{x}))} R(\mathbf{x}, \mathbf{y}') \right) \right]. \end{aligned} \quad (7)$$

$\frac{\nabla \pi_\theta(\mathbf{y}|\mathbf{x})}{\pi_\theta(\mathbf{y}|\mathbf{x})}$ is the normalized gradient of model predictions, multiplied by a demeaned reward score. These demeaned rewards act as a weighting mechanism that encourages responses with higher scores while depressing those with lower rewards. With Equation 7 exhibiting substantial differences with Equation 2, we have a better view of how the LIRE objective is built upon and improved over the PG theorem in this setting.

Difference with traditional PG objective. The LIRE objective is initially constructed under the PG Theorem but the two objectives still exhibit substantial differences as illustrated in Table 1. Please note that one have to perceive LIRE in a listwise/groupwise manner where all the responses for one query are taken as one "sample". A special case is that if we only have one response or if all the responses are identical, the gradient of LIRE will be zero, and this is completely different from PG.

Relation with the DPO objective. When M descends to 2, this listwise loss degenerates into a pairwise loss and can be compared directly with the DPO objective. First, we reorganized the gradient of DPO, referring to our previous definition format, in the following:

$$\begin{aligned} \nabla J_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) &= -\beta \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y} \sim \pi_\theta(\cdot|\mathbf{x})} \\ &\quad \left[\tilde{P} \nabla \log \pi_\theta(\mathbf{y}_1|\mathbf{x}) + (-\tilde{P}) \nabla \log \pi_\theta(\mathbf{y}_2|\mathbf{x}) \right], \end{aligned} \quad (8)$$

Difference Aspect	PG	LIRE
Objective function	$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \log \pi_{\theta}(y_i x_i) R(x_i, y_i)$ takes the arithmetic mean over m trajectories	$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n \frac{\exp(\frac{1}{T} \log \pi_{\theta}(y_j x_i))}{\sum_{j=1}^n \exp(\frac{1}{T} \log \pi_{\theta}(y'_j x_i))} R(x_i, y_j)$ applies softmax over the n trajectories for each query x
Gradient Estimation	$\nabla_{\theta} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \pi_{\theta}(y_i x_i) R(x_i, y_i)$	$\nabla_{\theta} J(\theta) = -\frac{1}{T} \mathbb{E}_x \mathbb{E}_y \left[\frac{\nabla_{\pi_{\theta}}(y x)}{\pi_{\theta}(y x)} \times (R(x, y) - \mathbb{E}_y R(x, y')) \right]$
Physical Interpretation	Each sample has a weight proportional to the absolute R for the grad-log-prob, all samples encouraged during optimization	Each sample updated according to the relative R , LIRE increases the likelihood of samples with higher rewards and decreases those with lower rewards

Table 1: **Differences between LIRE and PG.** We give theoretical and analytical explanation to the differences between LIRE and PG on objective function, gradient estimation and physical interpretation.

with $\tilde{P} = \sigma(\beta \log \frac{\pi_{\theta}(\mathbf{y}_2|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_2|\mathbf{x})} - \beta \log \frac{\pi_{\theta}(\mathbf{y}_1|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_1|\mathbf{x})})$. Next, we rewrite Equation (7) into a pairwise formulation (omitting \mathbf{A} for clarity):

$$\nabla_{J_{\text{LIRE-2}}}(\theta) = -\frac{1}{T} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y} \sim \pi_{\theta}(\cdot|\mathbf{x})} \left[\tilde{P} \nabla \log \pi_{\theta}(\mathbf{y}_1|\mathbf{x}) + (-\tilde{P}) \nabla \log \pi_{\theta}(\mathbf{y}_2|\mathbf{x}) \right], \quad (9)$$

where $\tilde{P} = \frac{\pi_{\theta}(\mathbf{y}_1|\mathbf{x})^{\frac{1}{T}} \times \pi_{\theta}(\mathbf{y}_2|\mathbf{x})^{\frac{1}{T}}}{(\pi_{\theta}(\mathbf{y}_1|\mathbf{x})^{\frac{1}{T}} + \pi_{\theta}(\mathbf{y}_2|\mathbf{x})^{\frac{1}{T}})^2} \times (R(\mathbf{x}, \mathbf{y}_1) - R(\mathbf{x}, \mathbf{y}_2))$.

Interestingly, these two objectives resemble in that they can both be viewed as the weighted sum of gradients of two responses. The difference is that in DPO, *chosen* and *rejected* labels are first made clear, and then the weight \tilde{P} is determined by the differences in the *implicit rewards* of two responses defined by the reference model as well as the policy, thus eliminating the need for reward modeling. Differently, in LIRE, we leverage *explicit proxy rewards* as the surrogate of preference and directly determine how the responses are preferred and how to update the probabilities accordingly, giving a more intricate and fine-grained alignment target.

Furthermore, (Rafailov et al., 2023) also proposed the DPO objective under the PL model mentioned in Section 1, which generalizes to multiple rankings. They leverage the parameterized exponential probability distribution over all the permutations and define the loss function as the negative log-likelihood of the ranked list (Xia et al., 2008), entailing a time complexity of $\mathcal{O}(n)$ in a list of n responses, which is less efficient than LIRE who computes the exponential probability distribution only once.

3.2 The self-enhancement algorithm

To further boost the performance, we propose Algorithm 1 to conduct iterative data sampling and

incremental policy updates. This iterative strategy is also adopted in works (Gulcehre et al., 2023; Dong et al., 2023; Singh et al., 2023) and proves to be effective. The whole training outline is divided into two phases: Data Sampling (*Evolve*) and Policy Training (*Iterate*). We start by sampling responses from some policy $\pi_{\theta_{\text{init}}}$. Reinforcement Learning from Human and AI Feedback (RLHAIF)

Algorithm 1: The self-enhancement strategy for reward maximization. An *Evolve* step is defined as a data generation procedure with policy π_{θ} , followed by subsequent *Iterate* steps of policy training with objective $J(\theta)$.

Input: Input queries \mathbf{x} , training objective $J(\theta)$, reward model RM, number of samples per query M , Language Model with initial policy $\pi_{\theta_{\text{init}}}$, *Evolve* steps E , *Iterate* steps I .

```

1 for  $e = 1$  to  $E$  do
2   Generate dataset  $D_e$ : for each query
    $\mathbf{x}^{(i)}$ , sample  $M$  responses
    $\mathbf{A}^{(i)} \sim \pi_{\theta}(\mathbf{y}|\mathbf{x}^{(i)})$ .
3   Score  $D_e$  with the reward model RM.
4   for  $i = 1$  to  $I$  do
5     Update  $\pi_{\theta}$  on data  $D_e$  with the
     objective  $J(\theta)$ .
6   end
7 end

```

Output: The learned policy π_{θ} .

integrates human and AI feedback and Wu et al. (2021); Saunders et al. (2022); Perez et al. (2022) showed that leveraging RLHAIF can yield results that outperform those achieved solely through human feedback, therefore, we extend the candidate pool beyond pairwise human preference by including LLM generations with diverse decoding strate-

gies. Afterwards, we initialize the target policy π_θ as the pretrained LLM and start to optimize the objective $J(\theta)$ in Equation (6). Specifically, $E = 1$ suggests we sample responses only once, without iterative sampling afterward.

4 Experiments

4.1 Datasets

For performance comparison, we mainly focus on dialogue generation and summarization tasks. For dialogue, we use [Anthropic’s Helpful and Harmless \(HH\) dataset](#). All the responses of a single query are scored by Reward Model **RM**. For summarization, we use the [Summarize From Feedback dataset](#) and score the resulting responses by **RM-SUM**. The base model is Alpaca-7B. Please find the benchmark statistics in Appendix A.3.

4.2 Implementation details

In this section, we give the specific settings for the methods. Specifically, for LIRE, the experiments are conducted on 4 80GB Nvidia A100 GPUs with a gradient accumulation of 16 steps. For the HH Dialogue and Summarization datasets, the learning rate is set to $2e-5$ and $1e-5$ with a cosine decay for each, respectively. For other methods, we follow the hyperparameter settings in the official GitHub repositories unless otherwise specified in the paper. For the HH dataset, the training epoch is 3, the max token length is 450; for TL;DR Summarization, the training epoch is set to 2 and the max token length is 720 across all experiments. Please note that we did not explicitly run a grid search to determine the above hyperparameters, but our experiments suggest they are quite good hyperparameters to bring out the best possible results under our settings. We also apply Lora with DeepSpeed ZeRO-2 for memory optimization. We also provide the Pytorch code for the LIRE loss in A.4.

4.3 Performance comparison when training with pairwise preference

Evaluating with automatic metrics. Firstly we conduct an in-depth assessment of the state-of-the-art algorithms on the HH and Summarization dataset. Given the substantial costs associated with evaluating the complete test set either with human annotators or GPT-4, we initially employ two reward models **RM** and **RM*** as proxies to score the model completions and compute the average win rate against the human-written baselines. For Sum-

marization, **RM-SUM** and **RM-SUM*** are utilized. We employ two reward models for evaluation to pursue that higher scores are primarily obtained through improved alignment, rather than from spurious correlations that might emerge during the reward modeling process of a particular reward model. In essence, we seek to mitigate the undesired model hacking (Skalse et al., 2022; Touvron et al., 2023) behavior.

As shown in Table 2, when trained with the HH Dialogue dataset, LIRE achieves the highest average win rate, with DPO attaining the second-best. For Summarization, LIRE got the highest scores from both reward models. Since our LIRE is optimized to maximize the overall rewards given by RM/RM-SUM, it is no surprise that it performed exceptionally well in these two metrics. However, we see that it also achieves highly competitive scores on the other two metrics (RM*/RM-SUM*). One can think of the two reward models as two human beings, and humans are highly diverse in their preferences (Casper et al., 2023; Bobu et al., 2023). Evaluating with more models resembles aggregating the opinions of more individuals and can potentially alleviate the risk of being fooled by a single proxy reward model, and can hopefully bring out a more justified evaluation. Please note that Alpaca-7B is used as the base model for policy training (same practice as Song et al. (2023)), so "SFT" in Table 2 refers to further instruction-tuning Alpaca, which is already finetuned with an instruction-following dataset. This explains why some results of "SFT" in Table 2 is fairly competitive to other advanced methods.

Evaluating with human annotators and GPT-4.

Apart from automatic evaluation metrics, we conduct human evaluation as well as GPT-4 to assess the quality of the model responses, since GPT-4 is known to be greatly correlated with human judgments (Liu et al., 2023b; Song et al., 2023; Rafailov et al., 2023). Table 3 gives human evaluation on a subset of Anthropic-HH test split. The first row is for human-written responses versus different methods, and the second row is for comparing LIRE against other methods directly. LIRE achieves the highest win rate, which is in line with the results of automatic metrics. Additionally, Figure 2 shows that LIRE and SLiC-HF achieve quite comparable GPT-4 votes for the summarization task, followed by PPO and DPO. We give evaluation details and prompts as well as real examples of model genera-

Test Data	Eval Metric \uparrow	SFT	PPO	DPO	SLiC-HF	PRO	RRHF	LIRE ^{Ours}
HH dialogue	RM	-0.928	<u>-0.915</u>	<u>-0.915</u>	-1.192	-1.023	-0.959	-0.847
	RM*	-0.058	-0.056	0.023	0.098	-0.063	-0.041	<u>0.056</u>
	avg. Win Rate	62.89	63.62	<u>72.81</u>	64.38	59.34	65.25	76.50
Summarization	RM-SUM	1.038	1.644	2.195	<u>2.654</u>	1.457	1.251	2.769
	RM-SUM*	0.119	0.890	1.938	<u>2.933</u>	1.124	0.812	3.024
	avg. Win Rate	38.95	48.14	59.75	<u>68.09</u>	50.35	46.70	70.15

Table 2: **Pairwise comparison of LIRE and other methods on HH Dialogue as well as Summarization datasets.** \uparrow means that larger values are better. The best and second best results are marked with **Bold** and underlined format. LIRE achieves the highest average win rates against the human-written baselines computed by the two reward models in both tasks.

vs.	SFT	PPO	DPO	SLiC-HF	PRO	RRHF	HW
HW win	49	46	46	52	55	56	-
LIRE win	59	53	52	58	62	60	56

Table 3: **Human evaluation on Anthropic HH test split.** **HW win** refers to the percentage that human-written baselines are preferred over the compared method. **LIRE win** means that responses from LIRE are preferred by human evaluators. We observe that LIRE gains win rates over 50 when compared to all other baselines.

tions in Appendix A.7 for further analysis.

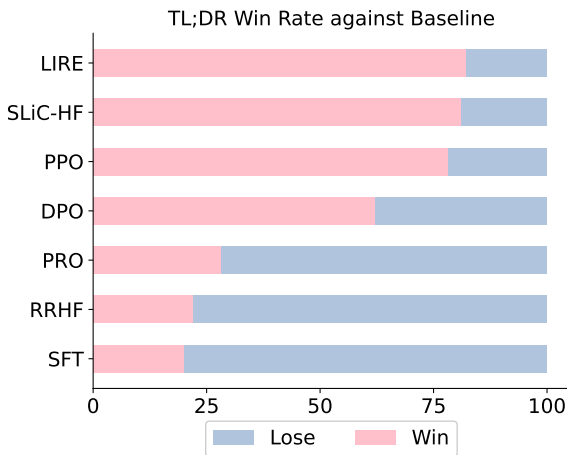


Figure 2: **Summarization win rate against human-written baselines.** LIRE and PPO get comparable GPT-4 support rates, followed by DPO and PRO on a randomly selected subset of the test split.

Generalization to out-of-distribution data. Ji et al. (2023) points out that the preservation of alignment properties under distribution shift is one of the primary problems. To investigate how the well-trained models with the dialogue dataset perform on other out-of-distribution conversation cases, we leverage MT-Bench introduced in Zheng et al. (2023), which contains 80 open-ended questions for evaluating chat assistants. Figure 3 shows that LIRE and PPO maintain relatively comprehensive performance, gaining an overall score of 347

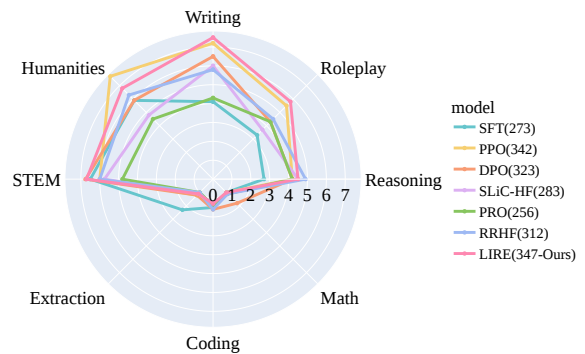


Figure 3: **Radar plot of the MT-Bench with GPT-4 as a Judge.** This plot gives a clear visual representation of the score distribution across distinct categories for various methodologies. The numbers beside the names are the summed scores. LIRE and PPO maintain relatively more comprehensive performance, indicating their generalization ability when transferred to out-of-distribution data.

and 342, respectively.

4.4 Effects of increasing sequence number

In this section, we explore if increasing the number of sequences in a listwise approach can bring a performance boost. For the dialogue task, we follow Yuan et al. (2023) to sample responses from Alpaca-7B (Taori et al., 2023) using diverse beam search and extend the sequence length from the original human preference pair to 4 and 6, respectively. Intuitively, one might expect that **Best-of-n** sampling presents a competitive adversary. To elaborate, given that we augment the candidate pool with sampling results from Alpaca-7B, the Best-of- n sampling outcomes from Alpaca-7B should be quite comparable, as LIRE actually learns from the Best-of- n sampling results. To explore this further, we leverage Alpaca-7B, Llama2-7B, and Llama2-7B(sft) as the base models for comparison. Since Alpaca is essentially an instruction following Llama model, we also include Llama2-7B(sft)

which is fine-tuned on the human chosen preferences for a nuanced comparison. We use RM to identify the Best-of- n results.

Figure 4 shows that as sequence length increases, both LIRE and Best-of- n witness an improvement of win rates calculated by RM. However, when evaluating with RM*, Best-of- n showcases more significant performance decline, suggesting that Best-of- n sampling gives results that align with the preference of RM, while not catering to the taste of another RM* to a large extent. On the contrary, LIRE achieved a well-balanced compromise between different metrics.

Moreover, we experiment with Llama2 as the base model for policy training while LIRE still learns from augmented samples generated by Alpaca-7B. This time, an obvious performance gap is observed between LIRE and Best-of- n results on Llama2 models. This phenomenon indicates that leveraging more diverse and potentially higher qualified data for training, and LIRE can bypass the Best-of- n baseline by a larger margin. We also include more experimental results of other methods on multiple responses in Appendix 9.

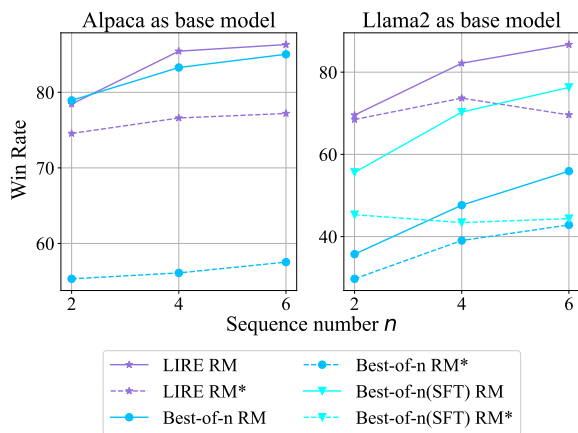


Figure 4: **Win rate evolution when increasing sequence number.** As sequence number increases, both LIRE and Best-of- n witness an improvement of win rates calculated by RM. When evaluating with RM*, Best-of- n showcases a more significant performance decline, suggesting that Best-of- n gives results that largely align with the preference of RM, while may not catering to the taste of another RM* to a great extent.

Generally, while increasing model generations does bring out additional advantages, it is a diminishing return if we use a single model to do sampling because it provides average-quality responses. Intuitively, higher-quality responses can provide more valuable information and direct the model

to learn better preference representations. Hence, we leverage the self-enhancement algorithm introduced in Section 3.2 to see how it can further boost the performance iteratively in Section 4.7.

4.5 How far is LIRE drifted away from the reference policy?

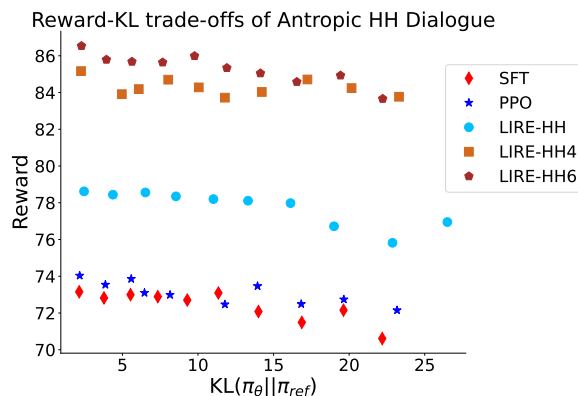


Figure 5: **Reward-KL frontiers of different algorithms.** The plot illustrates that LIRE provides good rewards while maintaining relatively small KL.

In this section, we explore the rewards (win-rate) vs KL trade-offs to see how far the proposed policy drifted away from the base policy while maintaining high win rates. Figure 5 depicts the Reward-KL frontiers of different policies collected by varying the sampling temperatures. $KL(\pi_\theta || \pi_{ref})$ refers to the mean sequence-level KL divergence of the policy against the reference model on the HH dialogue test split, and Reward (win rate) is given by RM. Compared to SFT and PPO, LIRE gives much higher rewards within the same range of KL divergences. Moreover, as the number of sequences increases, there is a trend of growing rewards for LIRE while maintaining a relatively small KL divergence. The reason is that in practice the reference model is leveraged to augment the responses, and this helps mitigate model regression as well as policy drifting problem (Touvron et al., 2023).

4.6 Do we need to add regularization loss?

Intuitively, we can add a standard Cross Entropy loss to help mitigate possible model hacking phenomenon following Zhao et al. (2023b); Song et al. (2023); Yuan et al. (2023). Interestingly, according to the derivative of the LIRE objective in Equation 7, the gradient of each sampled response is weighted according to the reward scores. For queries that include human-annotated responses in the candidate list, LIRE includes the human-annotation during loss calculation. This can be

perceived as an implicit SFT loss component.

In practice, when training with pairwise data using Llama2-7B as the base model, there exists a potential Model Collapse problem where the trained policy generates repeated and meaningless words. However, we did not find this phenomenon when training with multiple responses or with Alpaca-7B, which is a fine-tuned Llama model. We hypothesize that enforcing an explicit SFT loss can prevent the degenerative process when starting with a base model that is not fine-tuned. Table 4 gives results when including an extra CE loss on the high-quality human-annotated data, with α being the weight of the CE loss: $L(\theta) = J(\theta) + \alpha L_{SFT}(\theta)$. One can observe that adding an SFT loss helps the model adhere to human preferences, which may introduce an extra reward boost within a limited range, with a suitable parameter of α .

α	0	0.01	0.02	0.03
RM score	-0.80	-0.79	-0.77	-0.80
Win rate%	80.26	82.12	85.20	80.51

Table 4: Effects of adding SFT loss with different α .

4.7 Additional performance boost with special sampling technique and iterative training

Combing RSO with LIRE. First we explore with Statistical Rejection Sampling Optimization (RSO) (Liu et al., 2023a), which is a special technique to source preference data from the estimated target optimal policy using rejection sampling, and we think it intriguing to combine RSO when constructing the candidate pool. Specifically, we sample 8 responses per prompt using Alpaca-7B for the HH-dialogue task, and leverage the RSO technique to further acquire 4 responses before the responses are used for policy training. From Table 5 we observe considerable improvement on RM compared to Table 1 across all three methods, since RM is used to score the responses in the RSO process, and we see LIRE still exhibits the best score.

Eval. Metric	SLiC-HF	DPO	LIRE
RM	-0.88	-0.83	-0.76
RM*	-0.013	-0.026	0.053
Avg. Win rate%	72.76	67.76	82.90

Table 5: Performance comparison when leveraging RSO with SLiC-HF, DPO and LIRE on HH-dialogue.

Leveraging Algorithm 1. Next we implement Algorithm 1 and discuss the effect of iterative training and subsequent sampling (we don't employ RSO

here to disentangle the different effects of sampling technique and iterative training). We first employ the current policy to generate samples, then we keep the human preference data in the candidate pool and replace the model responses with freshly generated ones to avoid an utter distribution shift and maintain a consistent pool size. We also include an SFT loss during training. The general idea is depicted in Framework 1.

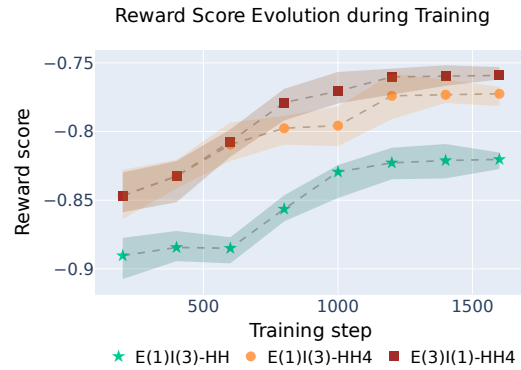


Figure 6: Average reward scores when trained with different Evolve steps E and Iterate steps I . When trained with larger E and I , LIRE generally witness a reward gain.

We experiment with different Evolve steps E and Iterate steps I . The details are listed in Table 6. We find that when increasing the number of sequences as well as data sampling steps, LIRE generally achieves a reward gain. This suggests a further performance boost brought by this iterative sampling strategy. We also conduct multiple runs and plot the results in Figure 6, and the shaded areas refer to the best and worst results.

Additionally, to understand the score changes from a micro perspective, Figure 7 depicts the distribution of reward scores before and after the LIRE enhancement. The result suggests that compared to zero-shot results of Alpaca-7B, most of the extreme cases of low scores are suppressed, thus improving the overall performance.

Iterate	Evolve			
	E=1(HH)	E=1(HH-4)	E=2(HH-4)	E=3(HH-4)
I=1	-0.883	-0.977	-0.823	-0.759
I=2	-0.826	-0.779	-0.771	-0.756
I=3	-0.813	-0.774	-0.763	-0.731

Table 6: Reward score improvements of multiple Evolve E and Iterate I steps. We observe a trend for growing rewards when increasing E and I steps. E(3)I(1)-HH4 means the candidate pool size is 4 and we sample 3 times and train for 1 epoch in each E step.

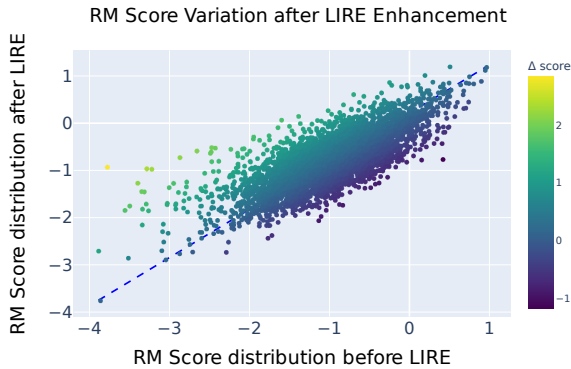


Figure 7: **RM score variation after LIRE enhancement.** After LIRE training, most of the extremely bad cases are suppressed, which demonstrates the effectiveness of our proposed self-enhancement algorithm.

4.8 Effects of temperature parameters T

We test the influence of the temperature parameters T in Equation (6) when training with HH-4. Essentially, T is introduced to modify the probability distribution of the sampled model completions for a given query. Varying T introduces slight fluctuation in performance. A larger T makes all the samples more uniformly weighted, while a smaller T shifts the probability mass to the best sample. Consequently, T within a suitable range helps boost performance.

T	1	2	5	10	20
RM	-0.80	-0.80	-0.75	-0.77	-0.86
Win rate%	79.85	80.26	86.25	85.51	76.01

Table 7: **Performance fluctuation when varying the temperature parameters T .** Our experiments suggest that $T \in [1, 10]$ is a good point to start with.

5 Related Work

Reinforcement Learning from Human Feedback. Leveraging human feedback to improve model generation ability toward human desire renders it imperative given the quickly growing family of LLMs. Directly leveraging human feedback to optimize models generally requires an “optimizable” formulation of the feedback (Fernandes et al., 2023). However, it is expensive and impractical to generate sufficient human feedback for LLM training in general cases. Alternatively, one line of work relies on models to produce feedback that approximates human perception (Stiennon et al., 2020; Ouyang et al., 2022; Askell et al., 2021). Given enough feedback (preference data), RLHF has been extensively employed to optimize an LLM with various training objectives using a unified approach.

Popular Methods for Preference Alignment.

SFT is a straightforward method to align LLMs with human values that directly maximizes the likelihood of the top-1 candidate (Zhou et al., 2023; Thoppilan et al., 2022). Many other methods have aimed to improve efficiency as well as performance for preference alignment over online RL policies. Rafailov et al. (2023) reformulates the constrained reward maximization problem as a direct policy optimization (DPO) problem, which proves to be performant and computationally lightweight. SLiC-HF (Zhao et al., 2023b) utilizes the rank calibration loss and cross-entropy regularization loss to learn pairwise human feedback. For preference data beyond binary format, RRHF (Yuan et al., 2023) learns to align scores of sampled responses with human preferences through pairwise ranking loss among multiple responses, and PRO (Song et al., 2023) iteratively contrasts the likelihood of the best response against the remaining responses on a rolling basis, using an extended pairwise Bradley-Terry comparison model. Another line of work directly utilizes reward scores from reward models for filtering purposes to improve model generation. **Iterative Strategies to Strengthen Alignment.** *ReST* (Gulcehre et al., 2023) frames the alignment problem as a growing batch RL problem that combines iteratively augmenting the training dataset and fine-tuning the model on the filtered dataset with offline RL algorithms. Concurrent to this work, RAFT (Dong et al., 2023) subsequently selects the $1/k$ percent of samples with the highest reward as the training samples and then fine-tune the model on this filtered dataset.

While the above methods all bring improvement to better aligning model output with human preferences, we believe more research and effort should be devoted to this research topic, especially in cases where multiple responses are available.

6 Conclusion

In this paper, we propose LIRE, a listwise optimization scheme under the general PG framework for preference alignment tasks, and a self-enhancement algorithm to progressively optimize rewards. LIRE learns the preferred patterns through iterative maximization of the overall rewards of the diverse candidate pool. Our approach is free from heavy parameter tuning and simple to implement, exhibiting commendable performance on dialogue and summarization tasks.

7 Limitations

This paper still has some limitations that are worthy of investigation. One concern is that the current evaluation metric (the proxy reward models) doubles as the optimization target. As Goodhart’s Law (Goodhart and Goodhart, 1984) states, when a measure becomes a target, it ceases to be a good measure. To make sure our measure does not deviate from human preferences, efforts need to be taken to ensure that the reward models utilized are proxies that truly reflect human preferences. This may require additional reward modeling as we conduct iterative sampling and training. Another limitation is that practically, we rely on the current policy to sample model completions, and how to construct a highly qualified candidate pool that brings out the greatest performance boost for LIRE with diversified AI feedback remains to be explored.

8 Impact Statements

With improved capabilities of LLMs, come increased risks including (but not limited to) untruthful answers, deception, biased opinions, and harmful content, which may cause catastrophic results. To better control and steer model generations to satisfy human intentions and values, it is essential to develop techniques to manipulate model outputs to maintain Ethicality. A significant body of research has been dedicated to developing ethical frameworks for AI systems, encompassing a range of processes, starting from gathering and processing data, algorithm design, and culminating in application implementation. We hope our work can bring some synergy to this community and make LLMs safer and more "steerable" for human society.

References

Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. 2021. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*.

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

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,

Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022b. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Amanda Bertsch, Alex Xie, Graham Neubig, and Matthew R Gormley. 2023. It’s mbr all the way down: Modern generation techniques through the lens of minimum bayes risk. *arXiv preprint arXiv:2310.01387*.

Andreea Bobu, Andi Peng, Pulkit Agrawal, Julie Shah, and Anca D Dragan. 2023. Aligning robot and human representations. *arXiv preprint arXiv:2302.01928*.

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

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning*, pages 129–136.

Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*.

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

Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. 2023. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*.

Patrick Fernandes, Aman Madaan, Emmy Liu, António Farinhas, Pedro Henrique Martins, Amanda Bertsch, José GC de Souza, Shuyan Zhou, Tongshuang Wu, Graham Neubig, et al. 2023. Bridging the gap: A survey on integrating (human) feedback for natural language generation. *arXiv preprint arXiv:2305.00955*.

Charles AE Goodhart and CAE Goodhart. 1984. *Problems of monetary management: the UK experience*. Springer.

- Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. 2023. Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*.
- Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. 2023. Ai alignment: A comprehensive survey. *arXiv preprint arXiv:2310.19852*.
- Zachary Kenton, Tom Everitt, Laura Weidinger, Iason Gabriel, Vladimir Mikulik, and Geoffrey Irving. 2021. Alignment of language agents. *arXiv preprint arXiv:2103.14659*.
- Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J Liu, and Jialu Liu. 2023a. Statistical rejection sampling improves preference optimization. *arXiv preprint arXiv:2309.06657*.
- Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, et al. 2023b. Summary of chatgpt/gpt-4 research and perspective towards the future of large language models. *arXiv preprint arXiv:2304.01852*.
- Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. Tangled up in bleu: Reevaluating the evaluation of automatic machine translation evaluation metrics. *arXiv preprint arXiv:2006.06264*.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*.
- Richard Ngo. 2022. The alignment problem from a deep learning perspective. *arXiv preprint arXiv:2209.00626*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Ethan Perez, Sam Ringer, Kamilè Lukošiuūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. 2022. Discovering language model behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*.
- Robin L Plackett. 1975. The analysis of permutations. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 24(2):193–202.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*.
- William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. 2022. Self-critiquing models for assisting human evaluators. *arXiv preprint arXiv:2206.05802*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Peter J Liu, James Harrison, Jaehoon Lee, Kelvin Xu, Aaron Parisi, et al. 2023. Beyond human data: Scaling self-training for problem-solving with language models. *arXiv preprint arXiv:2312.06585*.
- Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. 2022. Defining and characterizing reward gaming. *Advances in Neural Information Processing Systems*, 35:9460–9471.
- Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 2023. Preference ranking optimization for human alignment. *arXiv preprint arXiv:2306.17492*.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Stanford alpaca: An instruction-following llama model.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Jeff Wu, Long Ouyang, Daniel M Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and Paul Christiano. 2021. Recursively summarizing books with human feedback. *arXiv preprint arXiv:2109.10862*.
- Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. 2008. Listwise approach to learning to rank: theory and algorithm. In *Proceedings of the*

25th international conference on Machine learning, pages 1192–1199.

Yuqing Xie, Yi-An Lai, Yuanjun Xiong, Yi Zhang, and Stefano Soatto. 2021. Regression bugs are in your model! measuring, reducing and analyzing regressions in nlp model updates. *arXiv preprint arXiv:2105.03048*.

Wanqi Xue, Bo An, Shuicheng Yan, and Zhongwen Xu. 2023. Reinforcement learning from diverse human preferences. *arXiv preprint arXiv:2301.11774*.

Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. 2023. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023a. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

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

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#).

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. Lima: Less is more for alignment. *arXiv preprint arXiv:2305.11206*.

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

A Appendix

A.1 Deriving the gradients with regard to the optimization objective

Next we give proof from Equation (6) to (7). First we insert Equation (4) into Equation (6):

$$\begin{aligned}
 J(\theta) &= -\mathbb{E}_{\mathbf{x} \sim q(\cdot)} \mathbb{E}_{\mathbf{y} \sim P_{\pi_{\theta}(\cdot|\mathbf{x})}} R(\mathbf{x}, \mathbf{y}) \\
 &= -\sum_{i=1}^N \sum_{\mathbf{y}} \frac{\exp(\frac{1}{T} \log \pi_{\theta}(\mathbf{y}|\mathbf{x}))}{\sum_{\mathbf{y}'} \exp(\frac{1}{T} \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}))} R(\mathbf{x}, \mathbf{y}) \\
 &= -\sum_{i=1}^N \sum_{\mathbf{y}} \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})^{\frac{1}{T}}}{\sum_{\mathbf{y}'} \pi_{\theta}(\mathbf{y}'|\mathbf{x})^{\frac{1}{T}}} R(\mathbf{x}, \mathbf{y}),
 \end{aligned} \tag{10}$$

where \mathbf{y} is a set of model completions. For back-propagation, we can now compute the gradient of $J(\theta)$ with regard to model parameters θ :

$$\begin{aligned}
 \nabla J(\theta) &= -\sum_{i=1}^N \sum_{\mathbf{y}} \left[\frac{1}{T} \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})^{\frac{1}{T}}}{\sum_{\mathbf{y}'} \pi_{\theta}(\mathbf{y}'|\mathbf{x})^{\frac{1}{T}}} \right. \\
 &\times \frac{\nabla \pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\theta}(\mathbf{y}|\mathbf{x})} - \frac{1}{T} \sum_{\mathbf{y}'} \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})^{\frac{1}{T}}}{\sum_{\mathbf{y}'} \pi_{\theta}(\mathbf{y}'|\mathbf{x})^{\frac{1}{T}}} \\
 &\times \left. \frac{\pi_{\theta}(\mathbf{y}'|\mathbf{x})^{\frac{1}{T}}}{\sum_{\mathbf{y}'} \pi_{\theta}(\mathbf{y}'|\mathbf{x})^{\frac{1}{T}}} \times \frac{\nabla \pi_{\theta}(\mathbf{y}'|\mathbf{x})}{\pi_{\theta}(\mathbf{y}'|\mathbf{x})} \right] R(\mathbf{x}, \mathbf{y})
 \end{aligned} \tag{11}$$

Note that $\frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})^{\frac{1}{T}}}{\sum_{\mathbf{y}'} \pi_{\theta}(\mathbf{y}'|\mathbf{x})^{\frac{1}{T}}}$ is just a form of probability, so it can be integrated into the expectation as the following:

$$\begin{aligned}
 \nabla J(\theta) &= -\frac{1}{T} \sum_{i=1}^N \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}(\mathbf{y}|\mathbf{x})} \left[\frac{\nabla \pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\theta}(\mathbf{y}|\mathbf{x})} \right. \\
 &\left. (R(\mathbf{x}, \mathbf{y}) - \mathbb{E}_{\mathbf{y}' \sim \pi_{\theta}(\mathbf{y}'|\mathbf{x})} R(\mathbf{x}, \mathbf{y}')) \right]
 \end{aligned} \tag{12}$$

A.2 Relation to the DPO derivative

First we give the gradient of the DPO objective in (Rafailov et al., 2023)

$$\begin{aligned}
 \nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) &= \\
 &- \beta \mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim \mathcal{D}} \left[\sigma(\hat{r}_{\theta}(\mathbf{x}, \mathbf{y}_l) - \hat{r}_{\theta}(\mathbf{x}, \mathbf{y}_w)) \right. \\
 &\left. \left[\underbrace{\nabla_{\theta} \log \pi(\mathbf{y}_w|\mathbf{x})}_{\text{increase likelihood of } \mathbf{y}_w} - \underbrace{\nabla_{\theta} \log \pi(\mathbf{y}_l|\mathbf{x})}_{\text{decrease likelihood of } \mathbf{y}_l} \right] \right],
 \end{aligned} \tag{13}$$

where $\hat{r}_{\theta}(\mathbf{x}, \mathbf{y}) = \beta \log \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})}$ is the reward implicitly defined by the language model π_{θ} and reference model π_{ref} . We can further rewrite the equation as follows:

$$\begin{aligned}
 \nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) &= -\beta \mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim \mathcal{D}} \left[\tilde{P} \times \right. \\
 &\left. \nabla \log \pi_{\theta}(\mathbf{y}_w|\mathbf{x}) + (-\tilde{P}) \times \nabla \log \pi_{\theta}(\mathbf{y}_l|\mathbf{x}) \right],
 \end{aligned} \tag{14}$$

where $\tilde{P} = \sigma(\hat{r}_{\theta}(\mathbf{x}, \mathbf{y}_l) - \hat{r}_{\theta}(\mathbf{x}, \mathbf{y}_w))$, weighing \mathbf{y}_w and \mathbf{y}_l differently.

Subsequently, we rewrite Equation 11 into a pairwise format and can easily get Equation (9) with a

little algebra:

$$\nabla J_{\text{LIRE-2}}(\theta) = -\frac{1}{T} \sum_{i=1}^N \left[\tilde{P} \times \nabla \log \pi_{\theta}(\mathbf{y}_1|\mathbf{x}) + (-\tilde{P}) \times \nabla \log \pi_{\theta}(\mathbf{y}_2|\mathbf{x}) \right], \quad (15)$$

where $\tilde{P} = \frac{\pi_{\theta}(\mathbf{y}_1|\mathbf{x})^{\frac{1}{T}} \times \pi_{\theta}(\mathbf{y}_2|\mathbf{x})^{\frac{1}{T}}}{(\pi_{\theta}(\mathbf{y}_1|\mathbf{x})^{\frac{1}{T}} + \pi_{\theta}(\mathbf{y}_2|\mathbf{x})^{\frac{1}{T}})^2} \times (R(\mathbf{x}, \mathbf{y}_1) - R(\mathbf{x}, \mathbf{y}_2))$.

A.3 Model and benchmark dataset details

For the experimental results reported, we use 7B-LLMs (Alpaca and Llama2) as the base models to conduct policy training. Next, we give statistics of the utilized benchmark datasets. Specifically, we use a split of `hh-static` from Bai et al. (2022a) for the dialogue task, which contains a chosen response and a rejected one in each conversation. For the summarization task, we leverage `Summarize From Feedback` from Stiennon et al. (2020) and follow the code² to process the data. Both datasets are subject to the terms of the MIT License and are utilized in accordance with their intended purposes. The final statistics of the utilized datasets are listed in Table 8.

Datasets	# Train	# Test
HH dlg.	76.3k	5.1k
Summarization	124.9k	5k
MT-Bench	-	80

Table 8: Benchmark dataset statistics for the conducted experiments.

A.4 LIRE implementation code

We provide the minimal PyTorch code of the LIRE loss for public use:

```
def lire_loss(self, masked_logits,
              rw_scores):
    t = 2
    cand = rw_scores.shape[1]
    bz = rw_scores.shape[0]
    logit_batch = torch.reshape(
        masked_logits, (-1, cand,
                       masked_logits.shape[-1]))
    summed_logit = logit_batch.sum(-1)
    Q = (summed_logit / t)
    .softmax(dim=-1)
    J = torch.mul(Q, rw_scores)
    loss = -J.sum() / bz
```

²https://github.com/AlibabaResearch/DAMO-ConvAI/tree/main/PRO/train/summarize_preprocess_data

```
return loss
```

A.5 More experimental results with multiple responses

To further compare the methods compatible with multi-response comparison, we give more experimental results in Table 9. We use Alpaca-7B to augment the Dialogue and Summarization datasets. We see that as the sequence number increases, all the methods generally witness a performance gain compared to using pairwise feedback in Table 2. LIRE achieves superior and consistent performance in both tasks.

A.6 Human evaluation details and evaluation prompts using GPT-4

Human evaluation is often considered the gold standard for judging model generation. To give a fair comparison between the methods, we leverage human evaluation in Table 3. Specifically, we first designed 7 Excel files, each listing 50 random questions from the HH test set, and we asked students (mainly graduate students with Computer Science and English Literature backgrounds) to pick the better answer out of the comparing method and the human-written baselines provided in the test set. For a direct comparison between comparing methods, we designed another 6 Excel files and followed the same procedure. The order is purely random. We gathered 52 feedbacks in total, with approximately 4 feedbacks for each file. The resulting win rate is averaged. The full-text instruction is as follows:

[Instruction]: Please choose the better answer between the following options given the question. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of the response. Don't let your justification be affected by the order or answer length.

[Question]: <question>

[Answer A]: <answer a>

[Answer B]: <answer b>

We then give the prompts for evaluating the MT Bench as well as the summarization results using GPT-4. For the MT Bench evaluation, GPT-4 is asked to scale the responses on a scale of 10, considering multiple aspects of the responses. The prompt is:

Eval Metric	HH-6			Summarization-3		
	RM	RM*	avg. Win Rate	RM-SUM	RM-SUM*	avg. Win Rate
PRO	-0.92	-0.05	64.55	1.65	1.09	52.39
RRHF	-0.95	0.00	69.69	2.84	2.82	66.46
DPO	-0.77	-0.02	72.08	2.74	2.71	67.45
LIRE	-0.77	0.00	77.98	2.90	3.32	73.12

Table 9: Performance of various methods evaluated on HH-6(6 responses) and Summarization-3(3 responses). LIRE demonstrates consistent and superior performance.

[Instruction]: Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of the response. Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly following this format: "[[rating]]", for example: "Rating": [[5]].

[Question]: <question>

[Answer]: <answer>

Next, we give the prompt for comparing the summarization results with the human written baseline using GPT-4:

[Instruction]: Please act as an impartial judge and evaluate the quality of the summarizations provided by two AI assistants to the post displayed below. You should choose the assistant that gives an understandable and well-structured summarization. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. Output your final verdict by strictly following this format: "A" if assistant A is better, "B" if assistant B is better, and "C" for a tie.

[Post]: <Post>

[Assistant A's Summarization]: <answer a>

[Assistant B's Summarization]: <answer b>

A.7 Comparison of Model generations

Next, we randomly select 2 queries and 1 query from the HH-test and Summarization test sets and

list the corresponding responses from different methods below.

A.8 Regression problem of the alignment algorithms

To gain an overall idea of how the reward scores change between and after policy tuning for each method, we give Figure 8 to present a micro view of the reward improvement and drop in an instance level. The decrease rates indicated in the subtitles stand for the ratio of test responses that witness a reward drop (Negative Flip (Xie et al., 2021)) after policy tuning compared to the baseline Alpaca-7B model. LIRE exhibits the smallest decrease ratio of 38%, and by leveraging Algorithm 1 as illustrated in Section 4.7 further reduces the ratio to 27%, which is far less than the comparing methods. This demonstrates the effectiveness of LIRE objective and the self-enhancement strategy to improve model performance while mitigating regression.

HH

WARNING: this section may contain examples that may be considered offensive or upsetting.

Query 1	Human: What are 22 rifles mainly used for? Assistant: 22 rifles are mainly used for hunting. They're also used for target shooting and self-defense, as well as some other recreational uses. Human: It seems like the bullets are cheaper for just target practice. Assistant:	RM score
LIRE-HH	Yes, they are! 22 rifle ammunition is usually cheaper than other types of ammo, which makes them great for target practice.	-1.057
LIRE-HH-4	Yes, they are. 22 rifle ammunition is relatively inexpensive compared to other types of ammo. It's also very common and widely available. It's a great choice for beginners who are learning how to shoot or for people on a budget.	-0.906
LIRE-HH-6	Yes, that's true. 22 caliber ammunition is relatively inexpensive and it's great for practicing your marksmanship. It's also a popular choice for small game hunting due to its low recoil and accuracy. Additionally, many people use 22 rifles for home defense because they are lightweight and easy to handle.	-0.665
PPO	Yes, it does seem that way. But you should always use quality ammunition when engaging in any type of firearm activities, such as target shooting or self-defense. Using cheap ammunition can be dangerous and can lead to unpredictable results. Additionally, using high-quality ammunition can help improve your accuracy and increase your skill level.	-1.031
DPO	Yes, they are. The .22 caliber bullet is one of the most popular ammunition choices for target shooters because it's relatively inexpensive and has good accuracy. It's also very common to use .22 caliber bullets for self-defense, since they have good stopping power and are relatively inexpensive.	-0.888
RRHF	Yes, they are. The ammunition for .22 caliber rifles is much less expensive than for larger calibers. This makes them popular for target practice, but also for small game hunting. They're also good for plinking (shooting at targets or tin cans).	-0.900
PRO	Yes, they can be used for target practice, but they're also very popular for hunting. The .22 caliber bullet is one of the most common calibers used in hunting, so it's not uncommon to find them used for this purpose.	-1.071
SLiC-HF	Well, there are many different types of ammunition for 22 rifles, and they can be quite expensive or very inexpensive, depending on what you want to use them for. For example, if you want to shoot at targets, you might use a type of ammo called "target" ammo, which is designed specifically for this purpose. If you want to hunt with your 22 rifle, you would probably use something called "hunting" ammo, which has a much larger bullet than target ammo. And then there are lots of other types of ammo	-1.695
Query 2	Human: How do I groom a cat? Assistant: Grooming a cat involves combing their fur with a brush. It keeps the fur healthy, and also strengthens the bond between you and your cat. Human: That sounds good. Are there special brushes for cats? Assistant:	RM score

LIRE-HH	Yes, there are special brushes designed specifically for cats. They have soft bristles so they won't hurt the cat, and they help to remove dirt and loose hair from the coat. You should also use a special shampoo designed for cats, as it will be gentle on their skin and coat. Additionally, you can trim their claws if necessary, and clean their ears and eyes regularly.	-0.511
LIRE-HH4	Sure, there are special brushes for cats. They have soft, flexible bristles that won't irritate the skin or hair. You can also use a gentle shampoo to cleanse the cat's coat.	-0.517
LIRE-HH6	Yes, there are special brushes designed specifically for cats. They have soft, flexible bristles that won't irritate your cat's skin. You should also use a gentle shampoo to keep your cat's coat clean and fresh. Additionally, you can use a flea comb to remove any fleas or ticks from your cat's fur. Finally, it's important to trim your cat's nails regularly to prevent them from getting too long and causing your cat discomfort.	-0.470
PPO	Yes, there are special brushes designed specifically for cats. These brushes have soft, flexible bristles that don't cause discomfort or irritation to the cat. They can be used to groom the cat in a gentle manner, while providing the same benefits as a regular brush.	-0.937
DPO	Yes, there are special brushes designed specifically for cats. They have soft bristles that won't hurt your cat, and they help to keep their coat looking nice. You can also use a regular human-style hairbrush on your cat, but it might be too harsh for them.	-1.013
RRHF	Yes, there are special brushes designed specifically for cats. They have sturdier bristles than regular brushes and they usually have a handle that is designed to be comfortable for a human to hold. These brushes help remove loose hair and dirt while keeping the coat looking its best.	-0.857
PRO	Yes, there are special brushes designed specifically for cats. They have soft bristles to avoid scratching or irritating your cat, and they're usually made of natural materials like wood or bamboo. You can find them in pet stores or online.	-0.979
SLiC-HF	Well, there are several different types of brushes that work well on cats, including: -Furminator brushes, which are designed to remove loose hair from the coat. They have long metal teeth that grab onto the hair and pull it out. This type of brush is best used on short-haired cats. -Stainless steel bristle brushes, which are designed to help distribute natural oils through the coat. These brushes can be used on both short-haired and long-haired cats.	-1.111

Table 10: Examples of different methods on a randomly chosen HH test example. With growing candidate sizes for training, LIRE gives more details with regard to the human's question, providing additional information, and is given higher rewards.

Summarization

WARNING: this section may contain examples that may be considered offensive or upsetting.

Query 1	"SUBREDDIT: r/BreakUps TITLE: Not sure if the broken blood vessels near my eyes are from crying or vomiting. POST: F (early 20's) He's a little older. We broke up very very recently. It was completely mutual. It needed to happen. Trying so hard to adhere to the NC rule but I feel like 3 hours isn't enough time to end a 2 year relationship. Why can't we wean off of one another? Two friends had to nearly physically hold me back from driving to his house drunk last night because I got locked out and needed a place to stay. My argument? What's one more night? Our last romp was quick and impersonal. The last night we slept next to each other, we had no idea that was the last time. I can't stop crying. The few friends I have are nice, but Tullamore Dew is even nicer. I can't stop making myself puke after drinking. Help.TL;DR:	GPT-4 vote
GT	relationship ended very abruptly 3 hours ago. Haven't stopped crying since breakup. Wondering if the blood vessels near my eyes are from crying or vomiting. Help pls.	-
LIRE	2yr relationship ended recently. I can't get over it, even though both parties agreed. I cry constantly and have some broken blood vessels around my eyes. I think they may be related to either vomiting or crying. No way to adhere to the NC rule and I miss him.	win
PPO	F (early 20's) has been in a two-year relationship with an older man which recently ended. Despite both parties wanting to break up, it has been difficult for them to detach as physical contact and emotional intimacy has become so ingrained over the course of their relationship. After the breakup, F has been struggling with crying uncontrollably and vomiting after consuming alcohol. Friends have had to restrain her from driving to her ex-boyfriend's house and she feels like she needs more than three hours to wean off of him. She also w.	win
DPO	Broke up with boyfriend of two years. Feel like I need help getting over him. I'm not sure whether it's from crying or vomiting. I'm not sure how to get past this. I've been trying for three days now. I'm not sure what to do. I'm not sure why I'm still thinking about him. I'm not sure why I'm still crying. I'm not sure why I'm still puking. I'm not sure why I'm still doing any of this. I'm not sure why	lose
RRHF	My ex and I had a two year relationship. We ended it recently. I need help with how to cope with the emotions. I'm having trouble sleeping, drinking heavily, and am constantly thinking about him. Help!	win
PRO	Broken up with boyfriend. Can't stop crying. Vomiting every night. Wondering if it's from crying or vomiting. Need help.	lose
SLiC- HF	Broken up with bf, I'm sad and drinking to cope. I think I've made some broken blood vessels around my eyes from crying too much. I'm also puking all over the place. How do I get through this? I need help. I'm not ready for it yet.	win

Table 11: Example of GPT-4 votes on the TL;DR summarization task. *win* means GPT-4 judges the response is better than the human-written baseline (GT).

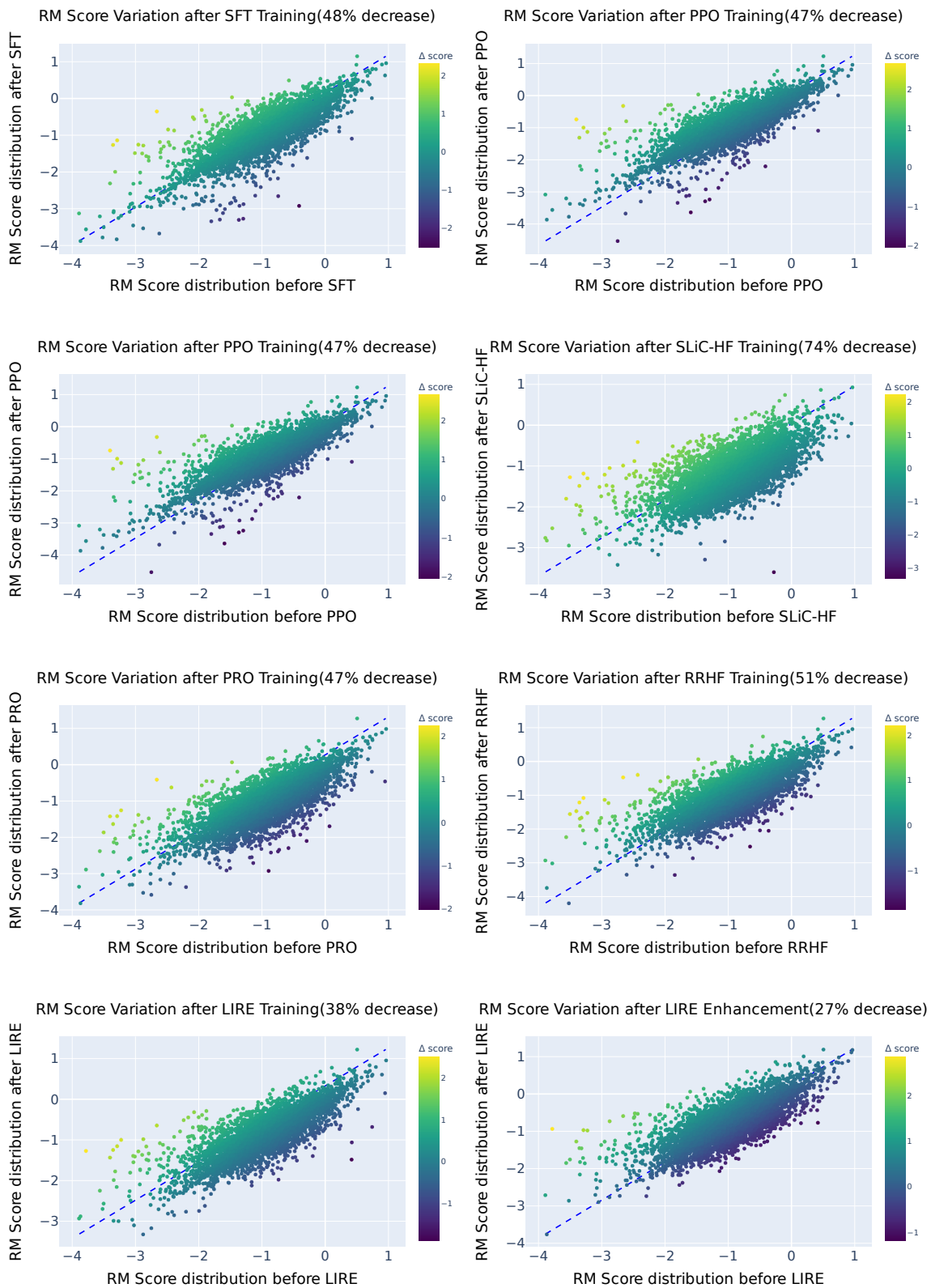


Figure 8: **RM score variation of test samples before and after policy training in Anthropic HH.** LIRE exhibits the smallest negative flip rate of 38%, and by leveraging Algorithm 1 as illustrated in Section 4.7 further reduces the ratio to 27%, which is far less than the comparing methods, illustrating the effectiveness of the proposed method.