

# Expectation Preference Optimization for Improving the Reasoning Capability of Large Language Models

Anonymous ACL submission

## Abstract

Pairwise preference optimization, such as Direct Preference Optimization (DPO), was originally designed to align large language models (LLMs) with human value. It has recently been used to improve the supervised fine-tuning (SFT) performance of LLMs. Using pairs of single samples, DPO estimates the probability distribution of the preferences of picking one response over another. However, in reasoning tasks that involve more complicated preferences than those in the human value alignment task, this sampling method is likely to bring deviations from the ground-truth distribution. To solve the problem, extra efforts (e.g., external annotations or amendment of the loss function) are often required. In this paper, we hypothesize that the preferences can be better estimated through a multi-sampling process. Accordingly, we propose an Expectation Preference Optimization (EPO) algorithm that takes pairs of sample groups, instead of pairs of single samples as in DPO, for preference learning. Compared to pairwise DPO, the proposed EPO tends to produce more proper preference estimations. Applying different preference optimization methods in a self-training paradigm, we have conducted extensive experiments on various reasoning benchmarks. The results show that our EPO approach outperforms a range of baseline approaches in terms of zero-shot accuracy on all benchmarks.

## 1 Introduction

Large language models (LLMs), through supervised fine-tuning (SFT), have shown remarkable abilities on various reasoning tasks such as mathematical reasoning. However, it is well recognized that the effectiveness of SFT can reach its upper limit depending on the scale and quality of training samples, which are often limited and expensive to construct. Thus an important question arises: *with the same SFT training data, how can we further*

*improve the SFT performance?* To tackle the problem, pairwise preference optimization, which was originally developed to align with human value (e.g. harmlessness or honesty), has become a widely chosen solution.

Direct Preference Optimization (DPO) (Rafailov et al., 2024) is one of the most popular preference-based methods due to its simplicity and effectiveness compared to Reinforcement Learning with Human Feedback (RLHF) (Bai et al., 2022). DPO samples the preferred and dis-preferred responses once in one updating step on a prompt, and then uses the Bradley-Terry (BT) model to update the LLM with an implicit reward function that models the preference of picking the preferred sample over the dis-preferred one. As can be naturally applied in the self-improving approaches that alleviate the issue of data construction (Yuan et al., 2024; Sun et al., 2023), using DPO in reasoning tasks has shown a broad prospect.

The selection of pairwise training data is key in the utilization of DPO. The preferred and dis-preferred responses on a prompt represent an estimation of the correct preference, which in the training process guides the optimization direction (Rafailov et al., 2024). Different from the human value alignment task, in most reasoning tasks the direction that the model needs to optimize can be more multifaceted. For example, in mathematical reasoning, the error of an answer can be attributed to various aspects, such as calculation, formula and entity errors. Thus directly using DPO on such reasoning tasks, especially when using correctness as the selection criterion for pairs of samples, would be insufficient to reflect the multifaceted nature of the reasoning tasks and result in a poor performance (Lu et al., 2024; Lai et al., 2024). As shown in Fig. 1 (the red box on the left-hand side), sampling a pair of single responses for optimization, with one reporting the correct answer and the other on the opposite, may lead to a

wrong direction of preference estimation that deviates from the other correct responses (marked with crying faces).

Various approaches have been developed to solve this problem. Orca-Math (Mitra et al., 2024) applies preference optimization on a fine-tuned LLM using an augmented dataset that is constructed using GPT4 to select the pairs of responses, while Brain (Chen et al., 2024a) uses human annotations. DPOP (Pal et al., 2024) tries to solve the unstable optimization direction of pairwise optimization by enhancing the supervision of preferred ends in changing the loss function of DPO. Iterative RPO (Yuanzhe Pang et al., 2024) uses a similar form of loss and applies it to a self-training structure. However, these methods do not fundamentally solve the problem of unstable preference modeling when facing complicated preferences.

In this paper, we explore a different perspective by *leveraging more samples in preference estimation*. Starting with the basic Bradley-Terry (BT) model, which is the basis of pairwise training, we hypothesize that the preferences in the BT model can be better estimated through a weighted multi-sampling process. Specifically, we assume that the preferences are not generated by the estimation of a single response, but by the expectation of the response sampling. Under this assumption, we propose an Expectation Preference Optimization (EPO) approach, a variant of DPO. EPO accepts group-wise preference samples, i.e., pairs of sample groups, for training, with a length limitation operation. EPO estimates the preference by calculating the weighted mean of each group. Our EPO shares the same objective with DPO and RLHF while overcoming the limitation of using only one preferred and one dis-preferred response each time. As shown in Fig. 1 (right-hand side), EPO makes it easier to produce proper preference estimations in reasoning tasks.

Utilizing the proposed EPO, we can simply use correctness as the signal for preference construction and boost the reasoning capability of LLMs yet bring no further human annotations. We apply a self-training algorithm which is detailed in Section 3.3. After SFT on a task-specific reasoning dataset, the target LLM generates responses for the input queries. Then we divide the responses for each query into two groups. Using EPO on these grouped responses, the optimization direction is estimated through multiple samples. Extensive experiments on various reasoning

benchmarks (i.e. GSM8K (Cobbe et al., 2021), ARC (Clark et al., 2018), SocialQA (Amini et al., 2019), MathQA (Sap et al., 2019)) across different base LLMs (including Llama2-7B, Llama2-13B (Touvron et al., 2023), Qwen1.5-7B (Bai et al., 2023), Mistral-7B (Jiang et al., 2023)) show that our EPO constantly improves the performance of SFT models and outperforms other preference optimization baselines in the self-training framework.

## 2 Preliminaries

Given a large language model that is parameterized by  $\theta$ , denoted as  $\pi_\theta$ , there are two categories of methods to improve its performance: fine-tuning-based and preference-optimization-based methods.

### 2.1 Fine-Tuning

**SFT:** Given a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ ,  $\pi_\theta$  is finetuned with the cross-entropy loss following a typical chain-of-thought rationale  $y_i$  with respect to the input query  $x_i$ , resulting in  $\pi_\theta^{SFT}$ .

**RFT:** Rejection Sampling Fine-Tuning (RFT) (Yuan et al., 2023) is a training method where  $\pi_\theta$  is fine-tuned on its own correct generations. After SFT on  $\mathcal{D}$ ,  $\pi_\theta^{SFT}$  obtains the ability to perform zero-shot chain-of-thought rationales. Thus we can sample  $M$  candidate rationales  $\hat{y}_{i,1}, \hat{y}_{i,2}, \dots, \hat{y}_{i,M}$  for each query  $x_i$ . All the rationales together are denoted as  $\hat{\mathcal{D}} = \{(x_i, \hat{y}_{i,j})_{j=1}^M \mid (x_i, y_i) \in \mathcal{D}\}$ .

Utilizing a filtering method (e.g. reward model annotation), we can construct  $\hat{\mathcal{D}}_{RFT}$  as a subset of  $\hat{\mathcal{D}}$ . The outcome  $\pi_\theta^{RFT}$  is trained on the augmented dataset  $\mathcal{D} \cup \hat{\mathcal{D}}_{RFT}$  based on  $\pi_\theta$ .

### 2.2 Preference-Optimization

**RLHF:** RLHF (Bai et al., 2022) fits a reward model to pairwise samples of human preferences and then uses Reinforcement Learning to optimize a language model policy to produce responses that are assigned high reward without drifting excessively far from the original model. Consider an annotated dataset of pairwise samples  $\mathcal{D}_p = \{x_i, y_w^i, y_l^i\}_{i=1}^N$ , where  $x_i$  denotes the  $i^{th}$  prompt,  $y_w^i$  and  $y_l^i$  respectively represent the preferred and dis-preferred responses to  $x_i$ . RLHF begins by modeling the probability of preferring  $y_w^i$  to  $y_l^i$  using the Bradley-Terry model (Bradley and Terry, 1952), which appoints the following probabilistic form:

$$p(y_w^i \succ y_l^i \mid x) = \sigma(r(x_i, y_w^i) - r(x, y_l^i)) \quad (1)$$

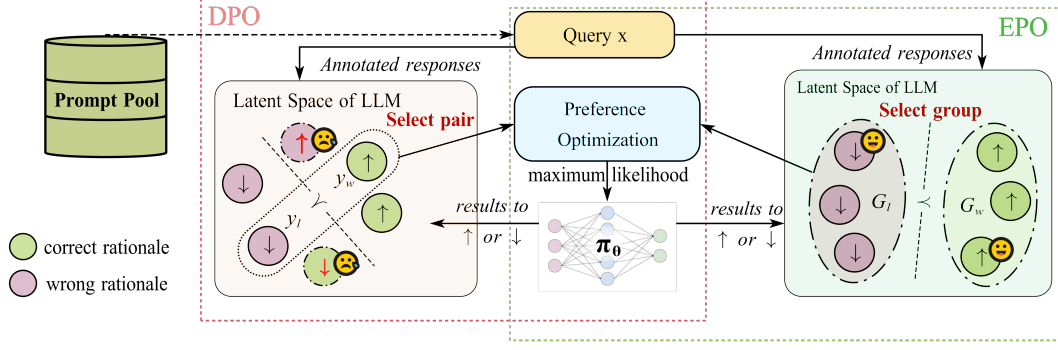


Figure 1: In the latent space of the target LLM, DPO chooses a pair of samples using correctness as the signal. In more complicated case, as shown in the figure, DPO can result in a wrong estimation of the preference and drive the LLM to a wrong reward updating direction (i.e., increased reward to the wrong samples and decreased to the correct samples). On the opposite, EPO considers multi-sampling and can provide a more reliable optimizing direction.

where  $\sigma$  represents the logistic function and  $r(x_i, y_i)$  corresponds to a reward function  $r_\phi$  (i.e., LLM classifier) that gives the estimation of  $y_i$  with respect to  $x_i$  according to human preference.

Then the target model  $\pi_\theta$  can be trained by the feedback from the learned reward function. In general, we formulate the following optimization target for this learning process:

$$\max_{\pi_\theta} \mathbb{E} [r_\phi(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi_\theta(y | x) \| \pi_{\text{ref}}(y | x)] \quad (2)$$

where  $\beta$  is a parameter controlling the deviation of the target model  $\pi_\theta$  from the status when the training starts.

**DPO:** DPO (Rafailov et al., 2024) shows the possibility of keeping the same optimization target as RLHF without explicitly training a reward function and the implementation of RL. The loss function of DPO is presented as below:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \log \sigma \left( \beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \quad (3)$$

Notably, this optimization objective is based on a theoretical optimal  $\pi_\theta$  beyond  $r_U(x, y)$ , which enables its equivalence with Eq.2.

### 3 Expectation Preference Optimization

#### 3.1 An Analysis of Pairwise Preference Optimization

Taking DPO as an example, Pairwise Preference Optimization methods accept one preferred sample

and one dis-preferred sample as the unit to calculate the loss for updating the reward function. Considering an ideal reward function  $\hat{r}(x, y)$  reflects the ground-truth preference, let us assume a sampling of four responses  $\{y_{\alpha 1}, y_{\alpha 2}, y_{\beta 1}, y_{\beta 2}\}$  with respect to the query  $x$ , where  $\hat{r}(x, y_{\alpha i}) > \hat{r}(x, y_{\beta i})$  holds. When an initial reward function  $r_\phi^t$  is optimized on  $(y_{\alpha 1}, y_{\beta 1})$ , the optimization directions of  $y_{\alpha 2}$  and  $y_{\beta 2}$  are not restricted to follow the ground-truth. The updated  $r_\phi^{t+1}$  may give a wrong estimation  $r_\phi^{t+1}(x, y_{\alpha 2}) < r_\phi^{t+1}(x, y_{\beta 2})$  while correctly estimating the training pair as  $r_\phi^{t+1}(x, y_{\alpha 1}) > r_\phi^{t+1}(x, y_{\beta 1})$ , and vice versa.

The trigger for this issue is that the sampling of  $(y_{\alpha 1}, y_{\beta 1})$  with respect to the prompt  $x$  may be away from the ground-truth preference distribution. Accordingly, the optimization of  $r_\phi^t$  gives wrong guidance on  $y_{\alpha 2}$  and  $y_{\beta 2}$ . When the purpose of training is to align with humans, the inconsistency of preference estimation is not so prominent (compared to reasoning tasks), so the problem is less significant. However, the reasoning tasks present a different situation. For example, in math reasoning tasks such as GSM8K, LLMs can make mistakes for many reasons (e.g., equation calculation errors, incorrect understanding of problems, etc.) and the estimates from different aspects are not independent. Thus the true preference distribution is complicated and varies with the target LLM.

#### 3.2 Expectation Preference Optimization

Aiming to solve the aforementioned problem brought by the single sampling of preference distribution in the reasoning tasks, we propose an Expectation Preference Optimization (EPO) algo-

rithm starting from the RLHF pipeline. As we have previously mentioned, the reward modelling phase of RLHF is based on the BT model. After a single sampling of response pair  $(y_1, y_2)$  for a prompt  $x$ , we can annotate the responses using human labellers or some stronger LLMs. As the preferences are presented as  $y_w \succ y_l \mid x$  where  $y_w, y_l \in \{y_1, y_2\}$  we can optimize a reward function through Eq. 1.

By estimating preferences through multi-sampling, which results in a group of responses  $\{y_{i=1}^N\}$  for a prompt  $x$ , we present the group-wise preference form  $G_w \succ G_l \mid x$  where  $G_w, G_l \subseteq y_{i=1}^N$ . In general,  $G_w$  represents the preferred group and  $G_l$  represents the dis-preferred group. We assume that the reward level of  $G_w$  and  $G_l$  is the expectation for all rewards in the group:

$$r^*(x, G) = \mathbb{E}_{y_i \sim G}[r(x, y_i)] \quad (4)$$

Thus the Bradley-Terry model can be rewritten as:

$$p^*(G_w \succ G_l \mid x) = \sigma(\mathbb{E}_{G_l}[r(x, y_i)] - \mathbb{E}_{G_w}[r(x, y_i)]) \quad (5)$$

**EPO objective.** Following the derivation process of DPO, we can construct the reward function under the optimal solution to Eq. 2 as follows:

$$r(x, y) = \beta \log \frac{\hat{\pi}(y \mid x)}{\pi_{\text{ref}}(y \mid x)} + \beta \log Z(x) \quad (6)$$

where  $Z(x) = \sum_y \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r(x, y)\right)$  represents a partial function referring to the previous work (Peters and Schaal, 2007; Rafailov et al., 2024). Using this re-parameterization of  $r(x, y)$ , Eq. 5 can be formed as below using the optimal solution.

$$p^*(G_w \succ G_l \mid x) = \sigma(\beta P_{G_l} - P_{G_w}) \quad (7)$$

$$P_G = \mathbb{E}_G[\log \frac{\pi(y_i \mid x)}{\pi_{\text{ref}}(y_i \mid x)}]$$

Due to space limitation, we present our proof and detailed deriving process in the Appendix.

We can now formulate a minimum loss function for the target model  $\pi_\theta$  through this preference function:

$$\mathcal{L}_R(r_\phi, \mathcal{D}) = -\mathbb{E}_{(x, G_w, G_l) \sim \mathcal{D}}[\log \sigma(P)] \quad (8)$$

While the sampling model (reference model) provides the group result (i.e.  $G_w, G_l$ ), we regard the  $\pi_{\text{ref}}(y_i \mid x)$  as the probability of  $y_i$  in the expectation. In practice, this means that the response with higher probability have a higher impact on the overall optimization direction. Thus, the loss function of EPO can be derived as:

$$\begin{aligned} \mathcal{L}_R(r_\phi, \mathcal{D}) &= -\mathbb{E}_{(x, G_w, G_l) \sim \mathcal{D}} \\ &[\log \sigma(\beta f(G_w, \pi, \pi_{\text{ref}}) - \beta f(G_l, \pi, \pi_{\text{ref}}))] \\ f(G, \pi, \pi_{\text{ref}}) &= \frac{\sum_{y_i \in G} \pi_{\text{ref}}(y_i \mid x)^\gamma \log \frac{\pi(y_i \mid x)}{\pi_{\text{ref}}(y_i \mid x)}}{\sum_{y_i \in G} \pi_{\text{ref}}(y_i \mid x)^\gamma} \end{aligned} \quad (9)$$

Notably, this method only calculates an approximate expectation, as the sum of probabilities is not 1. Thus we introduce a smoothing coefficient  $0 < \gamma \leq 1$ , to avoid weights with large variants caused by incomplete calculation of expected deviations.

**A further interpretation of EPO.** We here present a brief analysis of EPO. The objective function of EPO is derived from RLHF, which means that we share the same overall optimal solution with RLHF and DPO. As we estimate the preferences through a multi-sampling assumption, EPO has a more reliable implicit reward function compared to the pair-wise DPO, especially in reasoning tasks with complicated preferences. EPO drives the target LLM to have higher probabilities of generating responses in the preferred group and lower probabilities of generating responses in the dispreferred group, while ensuring the responses with higher probabilities affect more on the optimization. Notably, when the sampling number of  $G_l$  and  $G_w$  is 1, EPO becomes a typical DPO algorithm. Theoretically, in random sampling, the larger the sampling number, the more accurate the estimation of preferences in line with the ground-truth distribution.

**Length Limitation Operation.** After the brief analysis of the EPO loss function, we introduce an additional module to the EPO algorithm. Previous work (Wang and Zhou, 2024) indicates that the beginning tokens affect most of the decoding (generating) process of an LLM. Considering the subsequent tokens of the responses could adversely impact the coherence of the model in the optimizing process, especially the dis-preferred responses, we aim to increase the stability of the EPO optimization process by limiting the length of samples.

Specifically, we truncate the responses in  $G_l$  and  $G_w$  and ensure that the length of responses



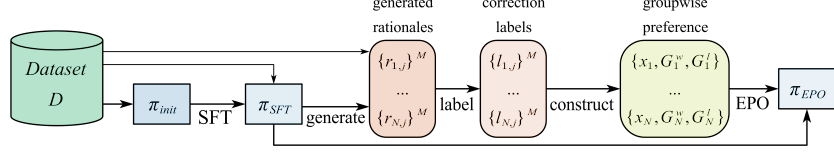


Figure 2: Overview of self-improving approach with EPO

is smaller than a preset threshold. Knowing that this truncation drops some information from the supervised data, we will analyze the effect of this operation in our experiments.

### 3.3 Self-improve Training approach With EPO

As EPO can provide reliable preference estimation, we can simply use correctness as signals and boost the reasoning capability of LLM. We design a self-improve training approach, which is presented in Fig 2.

We start with access to a base LLM  $\pi_{init}$  and samples of a reasoning task  $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$ . First, we give the model the ability to follow and generate rational instructions by applying SFT to it. The fine-tuned model is denoted as  $\pi_{SFT}$ . Then we generate  $M$  different responses for every query in  $\mathcal{D}$ . We denote all the generated responses ( $R_i$ ) with the original responses  $y_i$  as  $\mathcal{D}_{aug} = \{x_i, y_i, R_i\}_{i=1}^N$  where  $R_i = \{r_{i,j}\}_{j=1}^M$ .

In the next step, we generate the group-wise preference data from  $\mathcal{D}_{aug}$  using the correctness of generated responses in  $R_i$  as the annotation signal. Specifically, if a response reports the same answer as the typical rationale, it is put in  $G^w$  and it is put in  $G^l$  while it reports a different answer (which means it is wrong). The constructed training data are presented as follows:

$$\mathcal{D}_{EPO} = \{x_i, G_i^w, G_i^l\}_{i=1}^{N'} \quad (10)$$

where  $G_i^w \cup G_i^l = R_i \cup \{y_i\}$ . Notably, we construct the preference groups on  $R_i$  combining with  $y_i$ , thus for each prompt  $x$  the number of candidates' correct responses is always greater than 1. As the wrong response of a query does not always exist in the sampling, we drop the triplets in  $\mathcal{D}_{aug}$  whose  $R_i$  contains all correct responses.

Applying EPO algorithm on  $\pi_{SFT}$  with  $\mathcal{D}_{EPO}$ , we can obtain the resultant LLM denoted as  $\pi_{EPO}$ . In general,  $\pi_{EPO}$  is optimized based on the supervising information of base dataset  $\mathcal{D}$  (i.e. the correct answer), and the self-improve training en-

sures that the model can have better performance on the fine-tuning dataset.

## 4 Experiments

We evaluate the effectiveness of our EPO on two representative reasoning tasks: arithmetic reasoning and commonsense reasoning. We test four different LLMs: Llama2-7B (Touvron et al., 2023), Llama2-13B (Touvron et al., 2023), Qwen1.5-7B (Bai et al., 2023) and Mistral-7B (Jiang et al., 2023) as our base LLM model. We mainly evaluate the performance of EPO in the self-improving scenario. Notably, we put the our Implement Details in D

### 4.1 Datasets and Preprocessing

The experiments are carried out on two arithmetic reasoning datasets and three commonsense reasoning datasets.

**GSM8K.** GSM8K (Cobbe et al., 2021) has been adopted as a benchmark for the mathematical reasoning skills of LLMs. It contains 7,473 training and 1,319 test problems, and each sample is paired with a rationale that clearly states the final answer.

**MetaMath<sub>s</sub>** MetaMath (Yu et al., 2023) is a popular augmentation of GSM8K and MATH (Hendrycks et al., 2020). It contains 240K augmented samples based on GSM8K and 155K samples based on MATH. Notably, for lighter response generation, we only take 80K augmented GSM8K samples for training. The subset is denoted as MetaMath<sub>s</sub>.

**AI2 Reasoning Challenge (ARC).** ARC (Clark et al., 2018) consists of two subsets: ARC-Easy and ARC-Challenge. Each sample in the dataset contains a commonsense query and four candidate answers with one correct answer but does not contain rationales. To obtain the rationales of the queries for SFT, we apply a strong LLM (i.e. Yi-Chat-34B (Young et al., 2024)) to generate typical answers. Using the prompt presented in the Appendix, we generate a rationale ending with an answer statement for each query. After filtering the rationales with wrong answers and incorrect format,

we construct an SFT training set with 1599 samples from ARC-Easy and another with 793 samples from ARC-Challenge. They are then applied in the first SFT phase of the approach. For the generation phase, we use the original training set.

**MathQA.** MathQA (Amini et al., 2019) contains 29837 training samples and 2985 test samples. Each sample contains a math query, four candidate results, a rationale, and a correct answer. We manually add the answer statements at the end of the rationales for SFT.

**SocialQA.** Social IQA (Sap et al., 2019) has 33410 training samples, each containing a query and 3-5 candidate results without rationales, as well as 2224 test samples. We utilize the same method we use in constructing the ARC SFT dataset to generate rationales. Notably, we generate 23624 samples with one correct rationale each.

## 4.2 Baselines

In the experiments, we compare the proposed self-training EPO method (i.e. SFT + EPO) with various existing self-training approaches. They are described as follows. We present the detailed introduction in the Appendix.

**SFT** presents the  $\pi_{SFT}$  which is the LLM fine-tuned on typical rationales for specific tasks. It is used as the initialization of each self-training method below and our EPO.

(SFT +) **RFT** presents the model fine-tuned on the correct generated responses based on  $\pi_{SFT}$ , referring to the RFT method.

(SFT +) **DPO** presents the fine-tuned model using DPO on the pair-wise preference samples which are randomly chosen once for each prompt.

(SFT +) **DPO<sub>batch</sub>** presents the model using DPO training on pairs selected as many as possible to the prompt (while ensuring the single utilization of each response) in  $G_l$  and  $G_w$  for each prompt. It shows the performance of using batched DPO compared to EPO

(SFT +) **RPO** represents the model using the RPO algorithm (combining DPO loss with an NLL loss on the preferred response) on the pair-wise preference samples same as SFT + DPO.

## 4.3 Main Results

The main results of our experiments are presented in Tab. 1 and Tab. 2. Remarkably, on the GSM8K benchmark, EPO achieves a 5.43% increase over the SFT model in accuracy on the GSM8K dataset and 3.29% based on the Metasub<sub>s</sub> dataset for

Llama2-13B. This improvement comes to 4.54% and 7.05% for Qwen1.5-7B. As for the Commonsense tasks, EPO brings an increase of 6.01% for Llama2-7B on SocialQA, 4.47% for Mitral-7B on ARC-Easy, 6.94% for Llama2-13B on ARC-Challenge, and 6.29% for Mistral-7B on MathQA.

A cursory examination reveals that our EPO consistently outperforms all the preference optimization baselines across all tasks. Such a pattern underscores the effectiveness of EPO in improving LLM’s ability in reasoning tasks. The DPO baselines can eventually damage the performance of the model and this happens more frequently in mathematical reasoning. The DPO<sub>batch</sub> method also shows an unstable effect compared to the DPO while it can bring a slight improvement in many cases. RPO, compared to the former two, shows a more stable improvement effect. However, our EPO provides a more reliable preference estimation and constantly brings better performance improvements.

## 4.4 Further Analysis

### 4.4.1 Analysis of Generation Parameters and Length Limitation

**Effect of sampling temperature and length limitation.** We analyze the effect of sampling temperature in the generation phase and length limitation operation in the training phase. Fig. 3(a) shows the effectiveness of length limitation in contributing to the optimization stability. For GSM8K datasets, limiting the length of participation in the responses to the interval between 10 and 20 can result in better performance. As the sampling temperature grows, the peak is gradually moving rightwards. We consider this effect to be due to the increasing variety of responses that would decrease the instability responses.

**Effect of sampling number and length limitation.** We analyze the effect of sampling number in the generation phase and length limitation operation in the training phase. As shown in Fig 3(b), with the increase of the sampling number, the performance increases for the length limitation of less than 20. This result indicates that our EPO estimates the preference distribution more accurately as the number of samples increases. When the length limitation is increased, this benefit becomes unstable.

Table 1: Overall results on the math tasks in comparison with 4 base models. We report the accuracy of CoT Pass@1 greedy sampling. The best performance is in bold and the second-best is underlined.

Base Model	Datasets	SFT Result	Methods				
			RFT	DPO	DPO <sub>batch</sub>	RPO	EPO
Llama2-7B	GSM8K	28.96	<b>34.11</b>	28.45	28.47	27.89	<u>30.47</u>
	MetaMath <sub>s</sub>	<u>60.87</u>	60.27	59.34	58.65	58.21	<b>62.33</b>
Llama2-13B	GSM8K	49.27	47.99	48.47	48.53	<u>50.09</u>	<b>54.70</b>
	MetaMath <sub>s</sub>	69.82	68.38	67.39	68.46	<u>71.19</u>	<b>73.11</b>
Qwen1.5-7B	GSM8K	54.20	55.19	55.12	54.07	<u>56.46</u>	<b>58.74</b>
	MetaMath <sub>s</sub>	69.52	68.38	68.43	68.12	<u>70.56</u>	<b>76.57</b>
Mistral-7B	GSM8K	<u>41.84</u>	41.74	39.57	38.89	41.48	<b>45.40</b>
	MetaMath <sub>s</sub>	70.05	70.15	68.01	68.29	<u>71.72</u>	<b>74.72</b>

Table 2: Overall results on the Commonsense tasks in comparison with 4 base models. We report the accuracy of CoT Pass@1 greedy sampling. The best performance is in bold and the second-best is underlined.

Base Model	Datasets	SFT Result	Self-Training Methods				
			RFT	DPO	DPO <sub>batch</sub>	RPO	EPO
Llama2-7B	ARC-Easy	75.71	76.26	<u>77.23</u>	76.53	76.39	<b>78.10</b>
	ARC-Challenge	52.98	<b>56.56</b>	54.77	55.02	54.88	<u>55.74</u>
	MathQA	37.05	<u>38.15</u>	34.80	35.03	36.85	<b>38.88</b>
	SocialIQA	72.52	<u>72.52</u>	<u>77.42</u>	77.30	77.17	<b>78.53</b>
Llama2-13B	ARC-Easy	82.28	82.07	82.74	82.93	<u>83.20</u>	<b>84.35</b>
	ARC-Challenge	57.93	62.62	61.60	62.07	<u>63.99</u>	<b>64.87</b>
	MathQA	44.62	<b>47.07</b>	38.22	43.37	45.31	<u>46.91</u>
	SocialIQA	74.14	74.55	<u>78.50</u>	77.58	77.36	<b>79.86</b>
Qwen1.5-7B	ARC-Easy	85.35	85.85	87.74	87.03	<u>88.04</u>	<b>88.15</b>
	ARC-Challenge	74.82	70.32	77.55	<u>77.58</u>	76.73	<b>78.92</b>
	MathQA	51.89	51.75	51.62	<u>52.19</u>	50.84	<b>53.00</b>
	SocialIQA	73.74	74.82	75.85	<u>76.73</u>	<u>78.59</u>	<b>79.31</b>
Mistral-7B	ARC-Easy	74.47	72.83	74.83	75.05	<u>78.30</u>	<b>78.94</b>
	ARC-Challenge	60.45	62.71	63.84	60.03	<u>62.97</u>	<b>64.73</b>
	MathQA	52.09	52.36	50.83	51.95	<u>55.70</u>	<b>58.38</b>
	SocialIQA	74.10	74.37	<u>76.30</u>	75.58	76.15	<b>78.05</b>

#### 4.4.2 Effect of EPO from the Training Set Perspective

Considering that all the self-improving methods can more sufficiently utilize the training set compared to simple SFT, we analyze the performance of our EPO in comparison with baselines from the perspective of the training set. We apply an N=5 inference on GSM8K for each trained model with different methods. Taking the leftmost bar (SFT) in Fig 4 as the reference, we can observe that EPO increases the probability of the model responding correctly (i.e., increased number of the "5" segments and decreased number of the "0" segments) most. In fact, EPO drives the increase of the number of all-correct generations from 2441 to 3253, while DPO and RPO even drive it to decrease. This demonstrates the effectiveness of our method.

#### 4.4.3 Effects of Sampling Distribution on Training Result

As we utilize the expectation of a sampling process to estimate the preference in EPO, the sampling distribution (i.e. the samples in groups) can affect

Table 3: Effect of sampling distribution on DPO. "Highest / Lowest Prob" represents selection of the responses with the highest / lowest probabilities

Base Model	Random	Highest Prob	Lowest Prob
Llama2-7B	30.20	29.50(-0.70)	27.89(-2.31)
Llama2-13B	54.13	53.92(-0.21)	49.47(-4.66)

the final optimization direction. Here we present an analysis of the choice of responses for EPO. Firstly we apply an N=30 generation on GSM8K with T=0.7. Then we present three different methods to select 15 responses for each prompt: randomly selecting, selecting the responses with the highest probabilities, and selecting the responses with the lowest ones. We perform this analysis on two base LLMs: Llama2-7B and Llama2-13B. As shown in Table 4, the randomly selecting approach presents the best performance, and selecting with the lowest probabilities shows a poor performance. This implies that when selecting sample groups, it is necessary to follow a true distribution that guides a correct optimization direction, otherwise

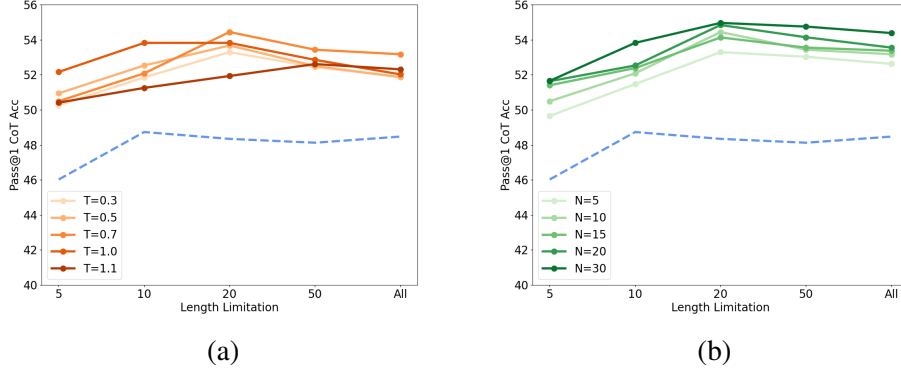


Figure 3: Analysis of hyperparameters. The analysis experiments are conducted on GSM8K for Llama2-13B. The sampling number for the experiments in (a) is set to 10 and the temperature for the experiments in (b) is set to 0.7. The blue dashed line represent the performance of DPO utilizing the length-limitation method.

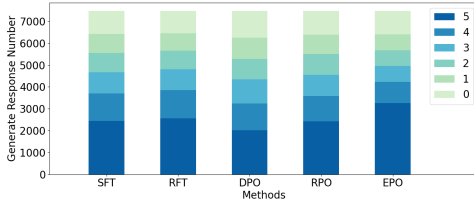


Figure 4: We calculate the number of correct responses for each query in an N=5 generation for each method on GSM8K, using Llama2-13B as base LLM. The different colors reflect different numbers of correct responses. The length of the bar represents the number of prompts.

optimization deviations may occur, leading to poor performance.

## 5 Related Work

Despite the success of instruction tuning on LLMs which has shown a great zero-shot performance (Chung et al., 2024; Mishra et al., 2021; Sanh et al., 2021), preference optimization has demonstrated its great effectiveness in aligning LLMs with humans (Bai et al., 2022). As reinforcement Learning with Human Feedback (RLHF) (Bai et al., 2022) is a complex and often unstable procedure (Pal et al., 2024), DPO (Rafailov et al., 2024) has been proposed as a more stable and computationally lightweight algorithm with no need for extra reward function training.

Reasoning ability is important for LLMs in practice. Let us take mathematical reasoning as an example. To make a stronger math-reasoning model, previous studies have focused on training the base model on larger datasets of better quality (Yuanzhe Pang et al., 2024; Yu et al., 2023). However, it is well-recognized that creating large-

scale and better-quality training samples is challenging and expensive.

The use of preference learning to improve the LLM’s reasoning ability has attracted increasing attention, while also facing certain problems. DPOP (Pal et al., 2024) enhances the supervision of the positive end in DPO by adjusting the loss function. Iterative RPO (Yuanzhe Pang et al., 2024) presents a similar loss function in a self-improving scenario without the SFT phase. Step-DPO (Lai et al., 2024; Lu et al., 2024) takes extra effort to create step-wise paired data and utilizes methods that are similar to vanilla DPO. However, these methods do not solve the problem of preference estimation of pair-wise optimization, thus gaining little improvement.

## 6 Conclusions and Future Work

In this paper, we propose an Expectation Preference Optimization (EPO) method that accepts pairs of response groups for preference learning. Compared to the existing pairwise preference optimization approaches that take pairs of single responses, our EPO method can more reliably estimate the preference distribution, especially when facing complicated reasoning tasks. We further design a self-improving framework, in which EPO can be effectively leveraged to improve the reasoning ability of LLMs. Experimental results on various reasoning tasks and datasets demonstrate the superior performance of our EPO which consistently outperforms a wide range of baseline approaches.

For future work, we plan to explore other reasonable methods (e.g., adding weights on responses) to better estimate the preferences based on EPO.



## 7 Limitations

Our paper presents a simple and practical method to improve the capability of LLMs in any reasoning task. However, the theory of EPO is not confined to reasoning tasks. Our intuition is to replace a single sample with an expectation in the Bradley-Terry model. Thus EPO can also be used in alignment tasks. However, we have not found a proper way to calculate the expectation in alignment tasks since in reasoning tasks the answer to a query is binary (i.e., correct or incorrect) while it is not in alignment tasks. Finding a proper method to calculate the expectation in alignment tasks can be a more comprehensive demonstration of the superiority of EPO theory.

## 8 Discussion of Ethical Considerations

Our proposed methods are used to improve the capabilities of LLMs. Though we mainly utilize it in reasoning tasks, it can also be used in other tasks which depends on the purpose of its user. On the other hand, using EPO training LLMs may cause an environmental impact as all other training methods do.

For the permissions of our used artifact, each of our used models (Llama2-13B, Llama2-7B, Mistral-7B, Qwen1.5-7B) and the datasets (GSM8K, ARC, MathQA) are open-sourced and can be found from Github or Huggingface. Secondly, all the models can not be used commercially.

We utilize all the models and datasets consistent with their intended use. We do not provide extra data. Our construction of self-training data using the LLMs presents the answers to the datasets, which is the purpose LLMs are designed.

The datasets we used contain no information that names or uniquely identifies individual people or offensive content.

## References

- Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*.
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. 2024. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- 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.
- Yezeng Chen, Zui Chen, and Yi Zhou. 2024a. Brain-inspired two-stage approach: Enhancing mathematical reasoning by imitating human thought processes. *arXiv preprint arXiv:2403.00800*.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024b. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind

690	Taffjord. 2018. Think you have solved question an-	Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xi-	744
691	swering? try arc, the ai2 reasoning challenge. <i>arXiv</i>	angru Peng, and Jiaya Jia. 2024. Step-dpo: Step-wise	745
692	<i>preprint arXiv:1803.05457</i> .	preference optimization for long-chain reasoning of	746
693	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	llms. <i>arXiv preprint arXiv:2406.18629</i> .	747
694	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias		
695	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri	748
696	Nakano, et al. 2021. Training verifiers to solve math	Edwards, Bowen Baker, Teddy Lee, Jan Leike,	749
697	word problems. <i>arXiv preprint arXiv:2110.14168</i> .	John Schulman, Ilya Sutskever, and Karl Cobbe.	750
698	Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff,	2023. Let’s verify step by step. <i>arXiv preprint</i>	751
699	Dan Jurafsky, and Douwe Kiela. 2024. Kto: Model	<i>arXiv:2305.20050</i> .	752
700	alignment as prospect theoretic optimization. <i>arXiv</i>		
701	<i>preprint arXiv:2402.01306</i> .	Zimu Lu, Aojun Zhou, Ke Wang, Houxing Ren,	753
702	Duanyu Feng, Bowen Qin, Chen Huang, Zheng Zhang,	Weikang Shi, Junting Pan, and Mingjie Zhan. 2024.	754
703	and Wenqiang Lei. 2024. Towards analyzing and	Step-controlled dpo: Leveraging stepwise error for	755
704	understanding the limitations of dpo: A theoretical	enhanced mathematical reasoning. <i>arXiv preprint</i>	756
705	perspective. <i>arXiv preprint arXiv:2404.04626</i> .	<i>arXiv:2407.00782</i> .	757
706	Alex Havrilla, Yuqing Du, Sharath Chandra Raparthy,	Yu Meng, Mengzhou Xia, and Danqi Chen.	758
707	Christoforos Nalmpantis, Jane Dwivedi-Yu, Maksym	2024. Simpo: Simple preference optimization	759
708	Zhuravinskyi, Eric Hambro, Sainbayar Sukhbaatar,	with a reference-free reward. <i>arXiv preprint</i>	760
709	and Roberta Raileanu. 2024. Teaching large lan-	<i>arXiv:2405.14734</i> .	761
710	guage models to reason with reinforcement learning.		
711	<i>arXiv preprint arXiv:2403.04642</i> .	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and	762
712	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	Hannaneh Hajishirzi. 2021. Cross-task generaliza-	763
713	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	tion via natural language crowdsourcing instructions.	764
714	2020. Measuring massive multitask language under-	<i>arXiv preprint arXiv:2104.08773</i> .	765
715	standing. <i>arXiv preprint arXiv:2009.03300</i> .		
716	Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron	Arindam Mitra, Hamed Khanpour, Corby Rosset, and	766
717	Courville, Alessandro Sordoni, and Rishabh Agar-	Ahmed Awadallah. 2024. Orca-math: Unlocking	767
718	wal. 2024. V-star: Training verifiers for self-taught	the potential of slms in grade school math. <i>arXiv</i>	768
719	reasoners. <i>arXiv preprint arXiv:2402.06457</i> .	<i>preprint arXiv:2402.14830</i> .	769
720	Hyeonbin Hwang, Doyoung Kim, Seungone Kim,	Arka Pal, Deep Karkhanis, Samuel Dooley, Man-	770
721	Seonghyeon Ye, and Minjoon Seo. 2024. Self-	ley Roberts, Siddhartha Naidu, and Colin White.	771
722	explore to avoid the pit: Improving the reasoning	2024. Smaug: Fixing failure modes of prefer-	772
723	capabilities of language models with fine-grained re-	ence optimisation with dpo-positive. <i>arXiv preprint</i>	773
724	wards. <i>arXiv preprint arXiv:2404.10346</i> .	<i>arXiv:2402.13228</i> .	774
725	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	Jan Peters and Stefan Schaal. 2007. Reinforcement	775
726	sch, Chris Bamford, Devendra Singh Chaplot, Diego	learning by reward-weighted regression for opera-	776
727	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	tional space control. In <i>Proceedings of the 24th in-</i>	777
728	laume Lample, Lucile Saulnier, et al. 2023. Mistral	<i>ternational conference on Machine learning</i> , pages	778
729	7b. <i>arXiv preprint arXiv:2310.06825</i> .	745–750.	779
730	Dahyun Kim, Yungi Kim, Wonho Song, Hyeonwoo	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	780
731	Kim, Yunsu Kim, Sanghoon Kim, and Chanjun Park.	Dario Amodei, Ilya Sutskever, et al. 2019. Language	781
732	2024. sdpo: Don’t use your data all at once. <i>arXiv</i>	models are unsupervised multitask learners. <i>OpenAI</i>	782
733	<i>preprint arXiv:2403.19270</i> .	<i>blog</i> , 1(8):9.	783
734	Diederik P. Kingma and Jimmy Ba. 2017. <a href="#">Adam:</a>	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	784
735	<a href="#">A method for stochastic optimization</a> . <i>Preprint</i> ,	pher D Manning, Stefano Ermon, and Chelsea Finn.	785
736	<i>arXiv:1412.6980</i> .	2024. Direct preference optimization: Your language	786
737	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	model is secretly a reward model. <i>Advances in Neu-</i>	787
738	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gon-	<i>ral Information Processing Systems</i> , 36.	788
739	zalez, Hao Zhang, and Ion Stoica. 2023. Efficient	Victor Sanh, Albert Webson, Colin Raffel, Stephen H	789
740	memory management for large language model serv-	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	790
741	ing with pagedattention. In <i>Proceedings of the 29th</i>	Chaffin, Arnaud Stiegler, Teven Le Scao, Arun	791
742	<i>Symposium on Operating Systems Principles</i> , pages	Raja, et al. 2021. Multitask prompted training en-	792
743	611–626.	ables zero-shot task generalization. <i>arXiv preprint</i>	793
		<i>arXiv:2110.08207</i> .	794
		Maarten Sap, Hannah Rashkin, Derek Chen, Ronan	795
		LeBras, and Yejin Choi. 2019. Socialliqa: Com-	796
		monsense reasoning about social interactions. <i>arXiv</i>	797
		<i>preprint arXiv:1904.09728</i> .	798

799	John Schulman, Filip Wolski, Prafulla Dhariwal,	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho,	853
800	Alec Radford, and Oleg Klimov. 2017. Proxi-	Sainbayar Sukhbaatar, Jing Xu, and Jason Weston.	854
801	mal policy optimization algorithms. <i>arXiv preprint</i>	2024. Self-rewarding language models. <i>arXiv</i>	855
802	<i>arXiv:1707.06347</i> .	<i>preprint arXiv:2401.10020</i> .	856
803	Zhiqing Sun, Yikang Shen, Hongxin Zhang, Qinhong	Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting	857
804	Zhou, Zhenfang Chen, David Cox, Yiming Yang, and	Dong, Chuanqi Tan, and Chang Zhou. 2023. Scal-	858
805	Chuang Gan. 2023. Salmon: Self-alignment with	ing relationship on learning mathematical reason-	859
806	principle-following reward models. <i>arXiv preprint</i>	ing with large language models. <i>arXiv preprint</i>	860
807	<i>arXiv:2310.05910</i> .	<i>arXiv:2308.01825</i> .	861
808	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho,	862
809	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	He He, Sainbayar Sukhbaatar, and Jason Weston.	863
810	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	2024. Iterative reasoning preference optimization.	864
811	Bhosale, et al. 2023. Llama 2: Open founda-	<i>arXiv e-prints</i> , pages arXiv-2404.	865
812	tion and fine-tuned chat models. <i>arXiv preprint</i>	Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B	866
813	<i>arXiv:2307.09288</i> .	Brown, Alec Radford, Dario Amodei, Paul Chris-	867
814	Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai	tiano, and Geoffrey Irving. 2019. Fine-tuning lan-	868
815	Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang Sui.	guage models from human preferences. <i>arXiv</i>	869
816	2023. Math-shepherd: Verify and reinforce llms	<i>preprint arXiv:1909.08593</i> .	870
817	step-by-step without human annotations. <i>CoRR</i> ,		
818	<i>abs/2312.08935</i> .		
819	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,		
820	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and		
821	Denny Zhou. 2022. Self-consistency improves chain		
822	of thought reasoning in language models. <i>arXiv</i>		
823	<i>preprint arXiv:2203.11171</i> .		
824	Xuezhi Wang and Denny Zhou. 2024. Chain-of-		
825	thought reasoning without prompting. <i>arXiv preprint</i>		
826	<i>arXiv:2402.10200</i> .		
827	Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yim-		
828	ing Yang, and Quanquan Gu. 2024. Self-play pref-		
829	erence optimization for language model alignment.		
830	<i>arXiv preprint arXiv:2405.00675</i> .		
831	Yuxi Xie, Anirudh Goyal, Wenye Zheng, Min-Yen		
832	Kan, Timothy P Lillicrap, Kenji Kawaguchi, and		
833	Michael Shieh. 2024. Monte carlo tree search boosts		
834	reasoning via iterative preference learning. <i>arXiv</i>		
835	<i>preprint arXiv:2405.00451</i> .		
836	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan,		
837	Lingfeng Shen, Benjamin Van Durme, Kenton Mur-		
838	ray, and Young Jin Kim. 2024. Contrastive prefer-		
839	ence optimization: Pushing the boundaries of llm		
840	performance in machine translation. <i>arXiv preprint</i>		
841	<i>arXiv:2401.08417</i> .		
842	Alex Young, Bei Chen, Chao Li, Chengen Huang,		
843	Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng		
844	Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi:		
845	Open foundation models by 01. ai. <i>arXiv preprint</i>		
846	<i>arXiv:2403.04652</i> .		
847	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu,		
848	Zhengying Liu, Yu Zhang, James T Kwok, Zhen-		
849	guo Li, Adrian Weller, and Weiyang Liu. 2023.		
850	Metamath: Bootstrap your own mathematical ques-		
851	tions for large language models. <i>arXiv preprint</i>		
852	<i>arXiv:2309.12284</i> .		

## 2 A Used Prompt

### 3 A.1 Prompt for Yi to generate rationales

4 user: Please answer the following single-choice question by presenting the thinking process and  
5 presenting the answer. 1. The question has an answer. 2. The thinking process part is a coherent  
6 paragraph. 3. Present the answer in the end of the response which is in the format of The answer is  
7 A/B/C/D.:

8 Question:

9 [present question here]

10 Choice:

11 [present choice here]

12 assistant:

### 13 A.2 Prompt for base models to generate CoT answer for GSM8K

14 Below is an instruction that describes a task.

15 "Write a response that appropriately completes the request.

16 Instruction:

17 [present query here]

18 Response:

### 19 A.3 Prompt for base models to generate CoT answer for Commonsense choosing task

20 Below is an instruction that describes a task.

21 Write a response that appropriately completes the request.

22 Instruction:

23 Pick the most correct option to answer the following question.

24 [present question here]

25 A.[present choice here]

26 B.[present choice here]

27 C.[present choice here]

28 D.[present choice here]

29 Response:

## 30 B Proof for optimal solution to EPO

### 31 B.1 Proof for optimal solution to EPO

32 We construct our proof following the previous works[1, 2]. From Eq. 2, our optimizing target is:

$$\max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi} [r(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi(y | x) \| \pi_{\text{ref}}(y | x)] \quad (1)$$

33 Notably, we can derive as:



$$\begin{aligned}
& \max_{\pi} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi} [r(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi(y | x) \| \pi_{\text{ref}}(y | x)] \\
&= \max_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[ r(x, y) - \beta \log \frac{\pi(y | x)}{\pi_{\text{ref}}(y | x)} \right] \\
&= \min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[ \log \frac{\pi(y | x)}{\pi_{\text{ref}}(y | x)} - \frac{1}{\beta} r(x, y) \right] \\
&= \min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[ \log \frac{\pi(y | x)}{\frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right)} - \log Z(x) \right]
\end{aligned} \tag{2}$$

where we define as :

$$Z(x) = \sum_y \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right) \tag{3}$$

Notably,  $Z(x)$  is a function of only  $x$  and  $\pi_{\text{ref}}$ . We can additionally define:

$$\pi^*(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right) \tag{4}$$

As is a probability distribution which holds  $\sum_y \pi^*(y | x) = 1$ . Using the  $Z(x)$ , we can re-organize the Eq. A1 as:

$$\begin{aligned}
& \min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \left[ \mathbb{E}_{y \sim \pi(y|x)} \left[ \log \frac{\pi(y | x)}{\pi^*(y | x)} \right] - \log Z(x) \right] = \\
& \min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} [\mathbb{D}_{\text{KL}} (\pi(y | x) \| \pi^*(y | x)) - \log Z(x)]
\end{aligned} \tag{5}$$

Since  $Z(x)$  does not depend on  $\pi$ , the optimal solution is achieved by the policy that minimizes the first term. The KL divergence is minimized in the situation where two distributions are equal. Thus we have the optimal solution:

$$\pi(y | x) = \pi^*(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right) \tag{6}$$

### B.1.1 Deriving the EPO Objective Under the Bradley-Terry Model

To derive the EPO objective under the Bradley-Terry preference model, we have the origin Bradley-Terry Model:

$$p^*(G_w \succ G_l | x) = \frac{1}{1 + \exp(\mathbb{E}_{y_i \sim G_l} [r(x, y_i)] - \mathbb{E}_{y_i \sim G_w} [r(x, y_i)])} \tag{7}$$

In Eq. 6, we have:

$$r(x, y) = \beta \log \frac{\hat{\pi}(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x) \tag{8}$$

Substituting Eq. A7 into Eq. A8, we can get:

$$\begin{aligned}
p^*(G_w \succ G_l \mid x) &= \frac{1}{1 + \exp(\mathbb{E}_{y_i \sim G_l} [r(x, y_i)] - \mathbb{E}_{y_i \sim G_w} [r(x, y_i)])} \\
&= \frac{1}{1 + \exp\left(\mathbb{E}_{y_i \sim G_l} \left[\beta \log \frac{\hat{\pi}(y_i|x)}{\pi_{\text{ref}}(y_i|x)} + \beta \log Z(x)\right] - \mathbb{E}_{y_i \sim G_w} \left[\beta \log \frac{\hat{\pi}(y_i|x)}{\pi_{\text{ref}}(y_i|x)} + \beta \log Z(x)\right]\right)} \\
&= \frac{1}{1 + \exp\left(\mathbb{E}_{y_i \sim G_l} \left[\beta \log \frac{\hat{\pi}(y_i|x)}{\pi_{\text{ref}}(y_i|x)}\right] - \mathbb{E}_{y_i \sim G_w} \left[\beta \log \frac{\hat{\pi}(y_i|x)}{\pi_{\text{ref}}(y_i|x)}\right]\right)} \\
&= \sigma\left(\mathbb{E}_{y_i \sim G_l} \left[\beta \log \frac{\hat{\pi}(y_i \mid x)}{\pi_{\text{ref}}(y_i \mid x)}\right] - \mathbb{E}_{y_i \sim G_w} \left[\beta \log \frac{\hat{\pi}(y_i \mid x)}{\pi_{\text{ref}}(y_i \mid x)}\right]\right)
\end{aligned} \tag{9}$$

46 This leads to Eq. 7.

## 47 C Implementation Details

### 48 C.1 Baselines

49 In this section, we present the details of the baselines we used compared to EPO. Notably, we are  
50 using different training methods in the self-training scenario. Thus all of our baselines start from the  
51 **SFT** model:

52 **SFT** presents the  $\pi_{SFT}$  which is the LLM fine-tuned on typical rationales for specific tasks. It is  
53 used as the initialization of each self-training method below and our EPO.

54 Beyond the **SFT** model, we utilize several self-training methods that do not introduce additional  
55 supervising information as our EPO does. The below methods are all beyond **SFT** model and the  
56 inference responses  $\hat{\mathcal{D}}$  sampled from **SFT** model and the certain dataset:

57 (SFT +) **RFT** presents the model fine-tuned on the correct generated responses based on  $\pi_{SFT}$ ,  
58 referring to the RFT method. Notably, we get a subset of  $\hat{\mathcal{D}}$  using the correction of responses as the  
59 filtering signal, denoted as  $\hat{\mathcal{D}}_{RFT}$ . **RFT** are fine-tuned on  $\mathcal{D} \cup \hat{\mathcal{D}}_{RFT}$ [3]. This method stands for the  
60 performance of fine-tuning in the self-improving scenario.

61 (SFT +) **DPO** presents the fine-tuned model using typical DPO on the pair-wise preference samples  
62 which are randomly chosen once for each prompt. Notably, we sample one correct response and one  
63 incorrect response for each prompt in  $\mathcal{D} \cup \hat{\mathcal{D}}$ [3] randomly. Then we apply DPO to this dataset. It has  
64 the same optimizing steps as our EPO.

65 (SFT +) **DPO<sub>batch</sub>** presents the model using DPO training on pairs selected as many as possible to  
66 the prompt (while ensuring the single utilization of each response) in  $G_l$  and  $G_w$  for each prompt.  
67 Notably, for each prompt in  $\mathcal{D} \cup \hat{\mathcal{D}}$ , we sample  $\min(Num_{right}, Num_{wrong})$  preference pairs as  
68  $Num_{right}$  and  $Num_{wrong}$  represent the number of correct and incorrect responses. It shows the  
69 performance of using batched DPO compared to EPO.

70 (SFT +) **RPO** represents the model using the RPO algorithm (combining DPO loss with an NLL loss  
71 on the preferred response) on the pair-wise preference samples same as SFT + DPO. Notably, the  
72 RPO objective is represented as:

$$\mathcal{L}_{RPO} = -\log \sigma \left( \beta \log \frac{M_\theta(c_i^w, y_i^w \mid x_i)}{M_t(c_i^w, y_i^w \mid x_i)} - \beta \log \frac{M_\theta(c_i^l, y_i^l \mid x_i)}{M_t(c_i^l, y_i^l \mid x_i)} \right) - \alpha \frac{\log M_\theta(c_i^w, y_i^w \mid x_i)}{|c_i^w| + |y_i^w|} \tag{10}$$

### 73 C.2 Hyperparameters

74 For the **SFT** training setups, we train SFT models using the following hyperparameters: learning rate  
75 of 2e-5, batch size of 64, max sequence length of 2048, and cosine learning rate schedule with 10%

warmup steps for 3 epochs. All the models are trained with an Adam optimizer [4]. This setting is also the same for **RFT**.

For the preference optimization **DPO**, **DPO<sub>batch</sub>**, **RPO** and **EPO**. We apply a search on the learning rate, training epoch, and additional hyperparameters. The search range is presented as below:

### C.3 Search range of Baselines

Table 1: Hyperparameter search range.

Methods	Search Range
<b>DPO</b>	$\beta \in [0.05, 0.1, 0.5, 1.0]$ $lr \in [1e-7, 2e-7, 5e-7, 1e-6]$
<b>DPO<sub>batch</sub></b>	$\beta \in [0.05, 0.1, 0.5, 1.0]$ $lr \in [1e-7, 2e-7, 5e-7, 1e-6]$
<b>RPO</b>	$\beta \in [0.05, 0.1, 0.5, 1.0]$ $lr \in [1e-7, 2e-7, 5e-7, 1e-6]$ $\alpha \in [0.25, 0.5, 1, 2]$
<b>EPO</b>	$\beta \in [0.05, 0.1, 0.5, 1.0]$ $lr \in [1e-7, 2e-7, 5e-7, 1e-6]$ $\gamma \in [0.1, 0.2, 0.5, 1.0]$

Notably, we are referring the papers [2, 5, 6] to set the search ranges. The length limitation of *EPO* is tuned from 5 to 100.

## D Implement Details

The experiments are carried out on 16 A100-80G GPUs with a Linux system. For all methods, we search the hyperparameters as we present the details in the Appendix. We train 3 epochs in each setting and report the performance of the best checkpoint. For the response generation phase in the self-improving scenario, we use the sample number  $N = 20$  with temperature  $T = 0.7$  following [5]. We use *Pytorch*<sup>1</sup> and *Huggingface*<sup>2</sup> as tools for the implementation. For preference optimization, we run our experiments based on *trl*<sup>3</sup>. All the generations were done using *vllm* [7]<sup>4</sup>. The code will be released on GitHub<sup>5</sup>.

## E The Time Cost of EPO

The training cost involves time cost and memory costs. For the former, taking the sample of 20 responses per prompt, EPO requires the LLM to process an input that is 10 times larger than other methods (20 to 2). Benefiting from CUDA’s parallel strategy for tensors, the extra time cost we need to bear is smaller than the linear estimation. For the latter, the extra GPU memory cost by a larger input tensor is much smaller than that is required for LLM training.

We present the relevance of training costs and the performance of our EPO. As it is shown in Fig 1, EPO’s training time is less than 3 times of the other methods (while  $N$  is less than 30), while requiring a small amount of extra GPU memory.

<sup>1</sup><https://pytorch.org/>

<sup>2</sup><https://huggingface.co/>

<sup>3</sup><https://github.com/huggingface/trl>

<sup>4</sup><https://github.com/vllm-project/vllm>

<sup>5</sup><http://github.com/xxxxxx>

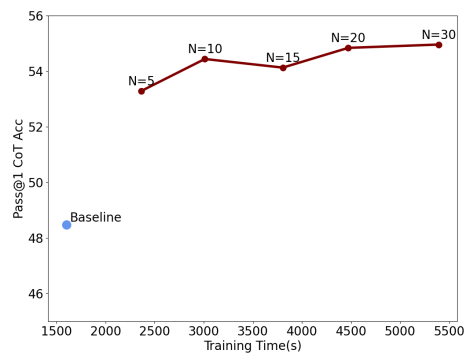


Figure 1: Analysis of training cost of EPO and baseline (i.e. DPO) under different N along with their performance.