

# Full-Step-DPO: Self-Supervised Preference Optimization with Step-wise Rewards for Mathematical Reasoning

Anonymous ACL submission

## Abstract

Direct Preference Optimization (DPO) often struggles with long-chain mathematical reasoning. Existing approaches, such as Step-DPO, typically improve this by focusing on the first erroneous step in the reasoning chain. However, they overlook all other steps and rely heavily on humans or GPT-4 to identify erroneous steps. To address these issues, we propose Full-Step-DPO, a novel DPO framework tailored for mathematical reasoning. Instead of optimizing only the first erroneous step, it leverages step-wise rewards from the entire reasoning chain. This is achieved by training a self-supervised process reward model, which automatically scores each step, providing rewards while avoiding reliance on external signals. Furthermore, we introduce a novel step-wise DPO loss, which dynamically updates gradients based on these step-wise rewards. This endows stronger reasoning capabilities to language models. Extensive evaluations on both in-domain and out-of-domain mathematical reasoning benchmarks across various base language models, demonstrate that Full-Step-DPO achieves superior performance compared to state-of-the-art baselines<sup>1</sup>.

## 1 Introduction

Large Language Models (LLMs) have attracted massive interest due to their remarkable capabilities across various tasks (Kaddour et al., 2023; Song et al., 2023; Wang et al., 2023a; Zheng et al., 2024; Wang et al., 2023b). However, they commonly encounter difficulties when tackling complex and symbolic multi-step reasoning, particularly in mathematical problem reasoning (Lightman et al., 2023; Huang et al., 2023). To improve the mathematical reasoning ability, some studies use Direct Preference Optimization (DPO) (Rafailov et al., 2024)

<sup>1</sup>Our code, data, and models are available at <https://github.com/anonymous>.

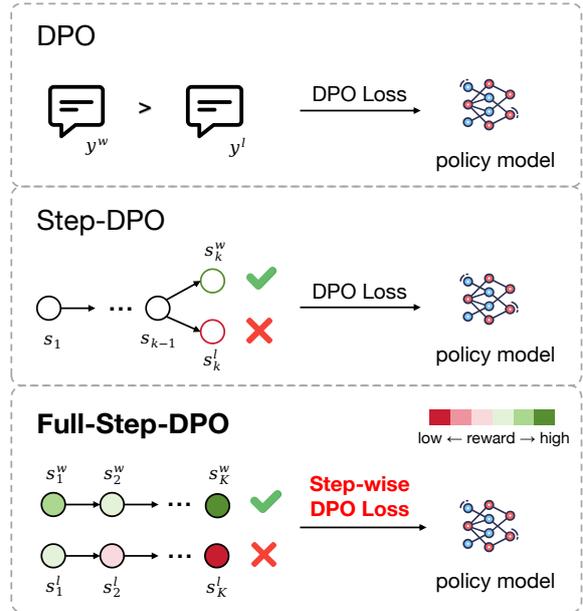


Figure 1: Comparison between DPO, Step-DPO, and our **Full-Step-DPO**. DPO operates on solution-wise preference data. Step-DPO advances to step-wise data but optimizes only a single step. Full-Step-DPO optimizes all steps with a novel step-wise DPO loss, effectively enhancing the model’s reasoning capability.

with solution-wise preference data but find its benefit limited (Pal et al., 2024; Xu et al., 2024; Jiao et al., 2024). Recent works attribute this limitation to DPO’s inability to perform process supervision and instead builds preference data based on reasoning steps rather than entire solutions (Lai et al., 2024; Chen et al., 2024; Xie et al., 2024a; Lu et al., 2024). For example, Step-DPO (Lai et al., 2024) focuses on optimizing only the first erroneous step in the reasoning chain, demonstrating notable improvements.

However, despite their improvements, these existing methods face the following limitations: (1) Some focus solely on the first erroneous step and ignore all other useful steps in the reasoning chain (Lai et al., 2024), as shown in Figure 1.

055 As a result, they fail to fully optimize the rea-  
056 soning chains, leading to suboptimal performance.  
057 (2) Their loss function still follows the early vanilla  
058 DPO in a solution-wise approach (Lu et al., 2024;  
059 Xie et al., 2024a). Consequently, it cannot directly  
060 leverage rewards in a step-wise fashion for learn-  
061 ing. (3) They heavily rely on costly and resource-  
062 intensive annotations from GPT-4 or humans to  
063 detect erroneous steps (Lai et al., 2024; Lightman  
064 et al., 2023), significantly limiting their practicality.

To address the above limitations, we propose  
065 **Full-Step-DPO**, a novel DPO framework for math-  
066 ematical reasoning. As illustrated in Figure 1, un-  
067 like vanilla DPO, which operates solution-wise, or  
068 Step-DPO, which focuses solely on the first erro-  
069 neous step, Full-Step-DPO utilizes each step in the  
070 entire reasoning chain and optimizes them using  
071 step-wise rewards. We first train a Process Reward  
072 Model (PRM) (Lightman et al., 2023; Wang et al.,  
073 2023c) in a self-supervised way, utilizing data gen-  
074 erated by the model itself. This approach enables  
075 the PRM to automatically score each step in the rea-  
076 soning chain, eliminating the reliance on external  
077 annotations such as GPT-4 or humans. Then, we  
078 propose a novel **Step-wise DPO Loss**, which em-  
079 ploys dynamic gradient updates to optimize each  
080 step based on its corresponding reward. This ap-  
081 proach shifts the optimization focus from solution-  
082 wise to step-wise, enabling the policy model to  
083 achieve superior reasoning capabilities.

We conduct experiments on both in-domain and  
084 out-of-domain mathematical reasoning datasets  
085 with four widely used backbone LLMs. Experi-  
086 mental results demonstrate that our Full-Step DPO  
087 consistently outperforms the DPO and Step-DPO  
088 baselines, validating its effectiveness in enhancing  
089 reasoning performance. Our contributions can be  
090 summarized as follows:

- 093 • We propose the Full-Step-DPO framework with  
094 a novel step-wise DPO loss that dynamically  
095 adjusts each step’s gradient based on its re-  
096 ward, enabling step-wise optimization rather than  
097 solution-wise and enhancing reasoning ability.
- 098 • We train a self-supervised PRM to provide step-  
099 wise rewards for preference learning and explore  
100 a more efficient approach for automatically con-  
101 structing PRM training data.
- 102 • Extensive experiments on widely used mathemat-  
103 ical benchmarks and base language models show-  
104 case the remarkable effectiveness of our method.

## 2 Related Work 105

**Mathematical Reasoning** Mathematical reason- 106  
ing task is one of the most challenging tasks for 107  
LLMs. Various approaches have been explored to 108  
improve or elicit the mathematical reasoning abil- 109  
ity of LLMs. A number of approaches have either 110  
continually pre-trained the base model on a vast 111  
of mathematical datasets (Azerbayev et al., 2023; 112  
Shao et al., 2024) or used supervised fine-tuning 113  
with substantial synthetic datasets distilled from 114  
cutting-edge models (Luo et al., 2023; Yu et al., 115  
2023b; Mitra et al., 2024; Xu et al., 2024). Another 116  
line of work focuses on enhancing test-time compu- 117  
tation by generating multiple solutions, developing 118  
separate reward models at either the outcome or 119  
process level to rerank these solutions (Cobbe et al., 120  
2021a; Lightman et al., 2023), or employing decod- 121  
ing strategies guided by the reward model (Yu et al., 122  
2023a; Xie et al., 2024b; Wang et al., 2023c; Wu 123  
et al., 2024). In addition, Reinforcement Learn- 124  
ing’s potential in general domains, demonstrated 125  
by Achiam et al. (2023) and Touvron et al. (2023), 126  
some studies have explored its use in mathematical 127  
reasoning (Wang et al., 2023c; Mitra et al., 2024; 128  
Pal et al., 2024). 129

**Preference Learning** Recently, preference learn- 130  
ing (Ethayarajh et al., 2024) has attracted signifi- 131  
cant attention due to its ability to align with human 132  
preferences and distinguish between positive and 133  
negative examples. While these methods, like DPO 134  
(Rafailov et al., 2024), have proven effective in gen- 135  
eral domains, it offers only marginal benefits for 136  
mathematical reasoning (Pal et al., 2024). Some 137  
works (Chen et al., 2024; Lai et al., 2024) suggest 138  
that DPO’s focus on coarse solution-wise prefer- 139  
ences makes it less effective at correcting errors 140  
in multi-step reasoning, hindering reasoning im- 141  
provement. Therefore, Step-DPO (Lai et al., 2024) 142  
was proposed, which first identifies the first erro- 143  
neous step, and then optimizes only this erroneous 144  
step along with the corresponding correct one. Al- 145  
though this approach enhances mathematical rea- 146  
soning capabilities, it totally overlooks the other 147  
steps in long-chain reasoning, which also provide 148  
valuable information and should not be completely 149  
disregarded. Building on this consideration, we 150  
propose Full-Step-DPO, which fully accounts for 151  
each step by dynamically optimizing all steps in 152  
the reasoning process. 153

**Step-wise Supervision** Recent findings by Lightman et al. (2023) suggest that step-wise supervision outperforms outcome-wise, due to the provision of more detailed feedback. However, training a PRM requires either costly manual annotation (Lightman et al., 2023) or significant computational resources (Khalifa et al., 2023; Wang et al., 2023c), which hinders the advancement and practical application of PRM. Therefore, in this paper, we aim to build a PRM for mathematical reasoning without relying on human annotation and with reduced computational resources. Additionally, we explore the effectiveness of the PRM in decoding and preference learning scenarios.

### 3 Full-Step DPO

In this section, we elaborate the proposed Full-Step DPO framework. We begin by reviewing the background of previous DPO and Step-DPO. Then we introduce the novel Step-wise DPO Loss which optimizes with step-wise rewards, and the Process Reward Model which automatically generate these step-wise rewards. Finally we outline the complete training pipeline of our Full-Step-DPO.

#### 3.1 Preliminary

**DPO.** Direct Preference Optimization (DPO) (Rafailov et al., 2024) is one of the most popular preference optimization methods. Instead of learning an explicit reward model, DPO directly uses pair-wise preference data to optimize the policy model with an equivalent optimization objective. Specifically, given an input prompt  $x$ , and a preference data pair  $(y^w, y^l)$ , DPO aims to maximize the probability of the entire preferred solution  $y^w$  and minimize that of the dispreferred solution  $y^l$ . The optimization objective of DPO is:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y^w | x)}{\pi_{\text{ref}}(y^w | x)} - \beta \log \frac{\pi_{\theta}(y^l | x)}{\pi_{\text{ref}}(y^l | x)} \right) \right]$$

where  $\pi_{\theta}(\cdot | x)$  is the policy model to be optimized,  $\pi_{\text{ref}}(\cdot | x)$  is the reference model,  $(x, y^w, y^l)$  are preference pairs,  $\sigma$  is the sigmoid function,  $\beta$  is a parameter controlling the deviation from the reference model.

**Step-DPO.** Although DPO performs well on chat benchmarks, it is less effective for long-chain reasoning tasks like mathematical problems. Step-DPO (Lai et al., 2024) attributes this to DPO’s

inability to consider the sequential nature of mathematical reasoning, as rejecting an entire dispreferred solution may inadvertently penalize correct preceding steps, introducing significant noise. To address this, Step-DPO optimizes only the first incorrect step. As shown in Figure 1, given a math problem and a series of initial correct reasoning steps  $\{s_1, \dots, s_{k-1}\}$ , Step-DPO aims to maximize the probability of the correct next step  $s_k^w$  and minimize the probability of the incorrect one  $s_k^l$ . Note that  $s_k^w$  and  $s_k^l$  refer to single steps, not all subsequent steps. The loss function used is still the vanilla DPO loss.

#### 3.2 Step-wise DPO Loss

We now introduce the novel Step-wise DPO loss, which performs step-wise optimization using step-wise rewards. Although the motivation behind Step-DPO is reasonable, focusing solely on optimizing the first erroneous step and neglecting the valuable information provided by other steps may not be optimal. Additionally, we contend that it is not truly a step-wise DPO, as it still relies on the standard solution-wise DPO loss and resembles more of a data construction method.

To address this, we modify the vanilla DPO loss to the step-wise DPO loss, dynamically weighting the gradients of each step based on its reward, thereby enabling true step-wise optimization. Let’s start with the gradient of the loss function  $\mathcal{L}_{\text{DPO}}$ . The gradient with respect to the parameters  $\theta$  can be written as:

$$\nabla_{\theta} \mathcal{L} = -\beta \mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}} \left[ \sigma \left( \hat{r}_{\theta}(x, y^l) - \hat{r}_{\theta}(x, y^w) \right) \left[ \nabla_{\theta} \log \pi_{\theta}(y^w | x) - \nabla_{\theta} \log \pi_{\theta}(y^l | x) \right] \right]$$

where  $\hat{r}_{\theta}(x, y) = \beta \log \frac{\pi_{\theta}(y | x)}{\pi_{\text{ref}}(y | x)}$ . Intuitively, the gradient indiscriminately increases the likelihood of whole  $y^w$  and decreases the likelihood of whole  $y^l$ . To achieve dynamically weighting, we break  $\nabla_{\theta} \log \pi_{\theta}(y | x)$  into a step-wise form and weight the gradient as follows:

$$\nabla_{\theta} \mathcal{L} = -\beta \mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}} \left[ \sigma \left( \hat{r}_{\theta}(x, y^l) - \hat{r}_{\theta}(x, y^w) \right) \left[ \sum_{i=1}^{K^w} \alpha_i^w \nabla_{\theta} \log \pi_{\theta}(s_i^w | x, s_{<i}^w) - \sum_{i=1}^{K^l} \alpha_i^l \nabla_{\theta} \log \pi_{\theta}(s_i^l | x, s_{<i}^l) \right] \right]$$

where  $s_i$  represents the  $i$ -th reasoning step of the solution  $y$ ,  $s_{<i}$  denotes all reasoning steps preceding  $s_i$ ,  $K$  is the total number of steps, and  $\alpha_i$  is

the weight coefficient of  $s_i$ , calculated based on the reward of  $s_i$  as shown below:

$$\alpha_i = \begin{cases} \frac{e^{\gamma r_{s_i}}}{\sum_j e^{\gamma r_{s_j}}}, & s_i \in y^w \\ \frac{e^{-\gamma r_{s_i}}}{\sum_j e^{-\gamma r_{s_j}}}, & s_i \in y^l \end{cases}$$

where  $r_{s_i}$  is the reward of the step  $s_i$ , which will be introduced in the next subsection, and  $\gamma$  is the temperature of the Softmax operation. It is important to note that the calculation of  $\alpha_i$  differs between the preferred solution  $y^w$  and the dispreferred solution  $y^l$ . For preferred solutions, a higher reward indicates a greater likelihood of correct reasoning in that step, so the model should perform gradient ascent with greater intensity. Conversely, for dispreferred solutions, a lower reward suggests a higher chance of incorrect reasoning, and the model should apply gradient descent with greater intensity accordingly. This approach allows us to leverage all steps and adaptively adjust the weight of each step based on its probability of correctness, achieving true step-wise optimization.

Compared to Step-DPO methods that focus solely on a single step, our method optimizes all steps simultaneously, enabling better global optimization. Noted that as  $\gamma \rightarrow 0$ , all steps will have equal weights, making Full-Step-DPO equivalent to vanilla DPO.

### 3.3 Process Reward Models

To obtain step-wise rewards, we train a Process Reward Model (PRM). The biggest challenge in training a PRM is constructing a process supervision dataset. Previous studies (Uesato et al., 2022; Lightman et al., 2023) utilize human annotators to obtain step-wise labels, which requires advanced annotator skills and is quite costly. Later, MathShepherd (Wang et al., 2023c) proposes using Monte Carlo estimation (Coulom, 2006) to automatically gather step-wise supervision, but it remained computationally expensive. In this section, we first examine the principles of Monte Carlo estimation, then present our simplified solution that significantly improves the efficiency of data construction.

**Monte Carlo estimation.** This approach assumes that the gold label  $y_{s_i}$  of a step  $s_i$  can be defined as the probability to deduce the correct answer  $a^*$ , and it includes both sampling and simulation phase. Specifically, given a math problem, it first

randomly samples  $M$  solutions, with each solution consisting of  $K$  reasoning steps  $\{s_1, s_2, \dots, s_K\}$ , and  $a$  represents the decoded answer from the last step  $s_K$ . Then, to estimate the quality of reasoning step  $s_i$  in a given solution, it simulates  $N$  subsequent reasoning processes from this step:  $\{(s_{i+1,j}, \dots, s_{K,j})\}_{j=1}^N$ . The golden label for  $s_i$  is calculated as follows:

$$y_{s_i} = \frac{\sum_{j=1}^N \mathbb{I}(a_j = a^*)}{N}$$

where  $a_j$  is the decoded answer for the  $j$ -th simulated solution, and  $\mathbb{I}$  is the indicator function that returns 1 if  $a_j = a^*$  and 0 otherwise. This two-stage approach is highly time-consuming, as it requires  $N$  simulations for each of  $K$  step across all  $M$  solutions, resulting in a time complexity of  $O(MNK)$ .

**Our efficient approach.** It is important to note that there is a trade-off between the sampling number  $M$  and the simulation number  $N$  when computational resources are limited. A larger  $M$  can provide more data for training the PRM, while a larger  $N$  can result in higher accuracy of the labels  $y_i$ . In this paper, we found that the trained PRM performs reasonably well even with  $N = 1$  when  $M$  is large, such as 32. This is likely because a larger  $M$  introduces more diversity into the training data, making the PRM more tolerant to slight reductions in data precision caused by the limited simulation number. This setting simplifies the PRM data construction by requiring only the sampling of  $M$  solutions without the need for simulation, significantly reducing computational resources and lowering the time complexity to  $O(M)$ . As a result, the gold label for step  $s_i$  can be simplified as follows:

$$y_{s_i} = \begin{cases} 1 & \text{if } a = a^* \\ 0 & \text{otherwise} \end{cases}$$

then the PRM could be trained as shown below:

$$\mathcal{L}_{\text{PRM}} = - \sum_{i=1}^K [y_{s_i} \log r_{s_i} + (1 - y_{s_i}) \log(1 - r_{s_i})]$$

where  $y_{s_i}$  is the golden label for  $s_i$ ,  $r_{s_i}$  is the sigmoid score assigned by the PRM. With the above PRM, we can automatically score each step in the reasoning chain, providing reward signals for the step-wise DPO loss and enabling step-wise optimization.

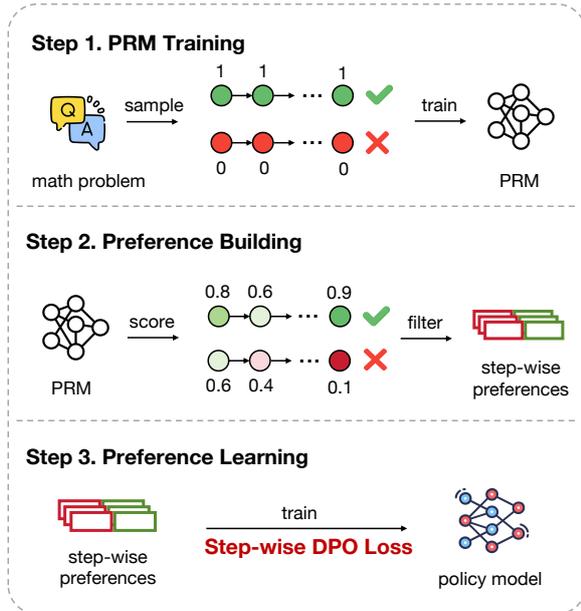


Figure 2: The overall framework of Full-Step-DPO consists of three steps: (1) Training the PRM using the model itself and generated solutions. (2) Using the PRM to score and filter solutions to form preference data with step-wise rewards. (3) Training the policy model with the proposed step-wise DPO loss.

### 3.4 Training Pipeline

Following previous methods (Wang et al., 2023c; Shao et al., 2024), we adopt a standard training pipeline illustrated in Figure 2: (1) We begin by training a PRM with self-generated data, where higher reward values indicate a stronger likelihood of correct reasoning, while lower values suggest potential errors. (2) The trained PRM is then used to construct preference pairs with step-wise rewards. Specifically, we generate  $M$  solutions for each math problem, score each step of these solutions with the PRM to produce a reward sequence, and calculate the average reward across all steps as the overall reward for each solution. We select the top  $T$  correct solutions with the highest rewards and the bottom  $T$  incorrect solutions with the lowest rewards to form  $T^2$  step-wise preference pairs. (3) Finally, we update the policy model using the proposed step-wise DPO loss and the step-wise preference pairs, as described in Section 3.2.

During the inference, a well-trained PRM can guide the decoding process and enhance the model’s performance. Therefore, in addition to the standard greedy decoding, we explore three alternative decoding methods: (1) Self-Consistency (SC) (Wang et al., 2022): given a problem in the test set, we sample  $K$  candidate solutions from the policy

model. Instead of relying on the first decoded solution, we select the final answer based on majority voting over the answers provided by all sampled solutions. SC is a simple yet highly effective verification strategy. (2) Best-of-N (BoN): we similarly sample  $K$  candidate solutions, score them using the reward model, and select the highest-scoring solution as the final answer. Following previous work (Lightman et al., 2023; Wang et al., 2023c), we use the minimum score across all steps as the final score assigned to a solution by the PRM. (3) Step-wise Beam Search (SBS) (Yu et al., 2023a): the PRM provides feedback at each step, offering more fine-grained guidance. Specifically, for each step, we first sample  $b_1$  candidate subsequent steps, then score them using the PRM. The top  $b_2$  steps are retained, and decoding continues until  $b_2$  final solutions are reached. The detailed algorithm is provided in Appendix A.

## 4 Experiments

### 4.1 Experimental Setup

**Backbones.** To comprehensively validate the effectiveness of our proposed method, we adopt four popular open-source LLMs as the backbone models: MetaMath-Mistral-7B (Yu et al., 2023b), Llama-3-8B (Touvron et al., 2023), DeepSeekMath-Base-7B (Shao et al., 2024) and Qwen2-7B (Bai et al., 2023). To improve these backbones’ reasoning ability, Step-DPO (Lai et al., 2024) finetunes DeepSeekMath-Base-7B and Qwen2-7B on two open-source synthetic math datasets, MetaMath (Yu et al., 2023b) and MMIQC (Liu and Yao, 2024), resulting in DeepSeekMath-Base-SFT<sup>2</sup> and Qwen2-7B-SFT<sup>3</sup>, which greatly outperform their previous versions. Following Step-DPO, we further finetune Llama3-8B to produce Llama3-8B-SFT. MetaMath-Mistral-7B has already been finetuned on MetaMath, so no additional finetuning was performed.

**Baselines.** For closed-source baselines, we compare our approach with OpenAI’s GPT-3.5 and GPT-4 (Achiam et al., 2023). We also benchmarked our method against recent high-performing mathematical LLMs, including WizardMath (Luo et al., 2023), MetaMath (Yu et al., 2023b), InternLM-Math-7B (Ying et al., 2024), Qwen2-7B-

<sup>2</sup><https://huggingface.co/xinlai/DeepSeekMath-Base-SFT>

<sup>3</sup><https://huggingface.co/xinlai/Qwen2-7B-SFT>

Instruct(Bai et al., 2023), DeepSeekMath-Instruct (Shao et al., 2024), InternLM-Math-20B (Ying et al., 2024), and Llama-3-70B-Instruct (Touvron et al., 2023).

Additionally, we compare it against DPO (Rafailov et al., 2024) and Step-DPO (Lai et al., 2024). Among these, Lai et al. (2024) publicly release DeepSeekMath-Base-SFT-Step-DPO<sup>4</sup> and Qwen2-7B-SFT-Step-DPO<sup>5</sup>, which we directly used for evaluation. Additionally, we trained MetaMath-Mistral-7B-Step-DPO and Llama-3-8B-SFT-Step-DPO using their publicly available code and dataset.

**Datasets.** To ensure a fair comparison, we use the same training dataset<sup>6</sup> provided by Step-DPO (Lai et al., 2024), which is synthesized from the training set of GSM8K (Cobbe et al., 2021b) and MATH (Hendrycks et al., 2021). Noted that we only use the problem prompts in this dataset and do not use the step labels marked by GPT-4.

For in-domain evaluation, we conduct experiments on GSM8K and MATH, which contain 1,319 and 5,000 test problems, respectively. We also evaluate on two more challenging out-of-domain (OOD) test sets OCWCourses (OCW) (Lewkowycz et al., 2022) and GaoKao2023 (GK2023) (Liao et al., 2024). OCW contains of 272 undergraduate-level STEM problems requiring multi-step reasoning for most questions, while GK2023 includes 385 mathematics problems from the 2023 Chinese higher education entrance exam, translated into English. Accuracy serves as the evaluation metric.

**Implementation Details.** During PRM training, we first randomly sample  $M = 32$  solutions for each math problem using Qwen2-7B-SFT and then label them as described in Section 3.3, resulting in the PRM training set. Then, we add a classification-head to Qwen2-7B-SFT and train it on the PRM training set for one epoch. The batch size is 256, and the learning rate is  $5e-7$ .

To build preference learning datasets, we first sample  $M = 32$  solutions for each math problem. The trained PRM then scores each solution, and

<sup>4</sup><https://huggingface.co/xinlai/DeepSeekMath-Base-SFT-Step-DPO>

<sup>5</sup><https://huggingface.co/xinlai/Qwen2-7B-SFT-Step-DPO>

<sup>6</sup><https://huggingface.co/datasets/xinlai/Math-Step-DPO-10K>

we select  $T = 4$  solutions with the highest average rewards and  $T = 4$  with the lowest average rewards to randomly form 16 preference pairs.

During preference learning, the batch size is 64, the learning rate is  $5e-7$ ,  $\beta$  is 0.05, and the reward temperature  $\gamma$  is 0.5. We use the AdamW (Loshchilov and Hutter, 2017) optimizer with a linear decay learning rate scheduler and only train one epoch. The warm-up ratio is 0.05.

During the decoding phase, we conduct experiments with two settings for SC and BoN, using  $K = 5$  and  $K = 15$ . For Step-wise Beam Search, to ensure fair comparison, we test two configurations:  $b_1 = 5, b_2 = 1$  (corresponding to  $K = 5$ ) and  $b_1 = 5, b_2 = 3$  (corresponding to  $K = 15$ ). The sampling temperature is set to 0.8.

All the experiments are conducted on a server equipped with 8 NVIDIA A100-80GB GPUs and 512GB of system RAM. The implementation frameworks are PyTorch (Paszke et al., 2017), DeepSpeed (Rasley et al., 2020), and Huggingface (Wolf et al., 2019).

## 4.2 Main Results

Table 1 provides a comprehensive comparison of various models on both MATH and GSM8K, including open-source and closed-source LLMs. We find that: (1) Consistent with previous studies (Pal et al., 2024), DPO exhibits notable instability. Its performance shows slight degradation on MetaMath-Mistral-7B and MetaMath-Mistral-7B-SFT backbones, while the accuracy drops sharply to around 20% on Qwen2-7B-SFT. It achieves a slight performance improvement only when applied to the DeepSeekMath-Base-SFT. (2) Step-DPO achieves only minimal improvements across all backbones, with gains generally around 1% and, in some settings, even slight performance drops. We evaluate the publicly released Step-DPO model using its official script, and the results may differ slightly from those reported in the Step-DPO paper. Similar issues have also been observed by other researchers<sup>7</sup>. (3) Our Full-Step-DPO consistently outperforms Step-DPO across all backbones. Specifically, when applied to MetaMath-Mistral-7B and Llama-3-8B-SFT, our model achieves improvements of approximately 2.3% to 3.7%, while applied to the stronger backbones, DeepSeekMath-Base-SFT and Qwen2-7B-SFT, our method still delivers gains exceeding 1%. These results clearly

<sup>7</sup><https://github.com/dvlab-research/Step-DPO/issues/2>

Model	MATH (%)	GSM8K (%)
GPT-3.5	34.1	80.8
GPT-4	53.6	93.6
WizardMath	10.7	54.9
MetaMath	19.8	66.5
InternLM-Math-7B	34.6	78.1
Qwen2-7B-Instruct	49.6	82.3
DeepSeekMath-Instruct	46.8	82.9
InternLM-Math-20B	37.7	82.6
Llama-3-70B-Instruct	50.4	93.0
MetaMath-Mistral-7B	28.2	77.7
+ DPO	24.8 <sup>-3.4</sup>	70.7 <sup>-7.0</sup>
+ Step-DPO	28.9 <sup>+0.7</sup>	79.6 <sup>+1.9</sup>
+ Full-Step-DPO	<b>30.5</b> <sup>+2.3</sup>	<b>81.4</b> <sup>+3.7</sup>
Llama-3-8B-SFT	32.6	78.5
+ DPO	23.4 <sup>-9.2</sup>	62.3 <sup>-16.2</sup>
+ Step-DPO	31.8 <sup>-0.8</sup>	80.1 <sup>+1.6</sup>
+ Full-Step-DPO	<b>35.0</b> <sup>+2.4</sup>	<b>82.0</b> <sup>+3.5</sup>
DeepSeekMath-Base-SFT	51.7	86.4
+ DPO	51.7 <sup>-0</sup>	87.3 <sup>+0.9</sup>
+ Step-DPO	52.9 <sup>+1.2</sup>	86.6 <sup>+0.2</sup>
+ Full-Step-DPO	<b>53.2</b> <sup>+1.5</sup>	<b>87.9</b> <sup>+1.5</sup>
Qwen2-7B-SFT	53.9	88.3
+ DPO	20.0 <sup>-23.9</sup>	27.3 <sup>-61.0</sup>
+ Step-DPO	54.9 <sup>+1.0</sup>	88.4 <sup>+0.1</sup>
+ Full-Step-DPO	<b>55.4</b> <sup>+1.5</sup>	<b>89.3</b> <sup>+1.0</sup>

Table 1: Performance comparison of various models on MATH and GSM8K with greedy decoding.

demonstrate the effectiveness of our proposed approach, which considers all steps in the reasoning process rather than focusing on solution-wise preferences or only a single step. A case study can be found in Appendix B.3.

### 4.3 Results on OOD Datasets

To further demonstrate the superiority of Full-Step-DPO, we evaluate the models on OOD datasets GK2023 and OCW, as shown in Table 2. On these competition-level math problems, DPO and Step-DPO often exhibit performance degradation under various settings, while our Full-Step-DPO consistently achieves performance improvements. The only exception occurs on the OCW dataset with MetaMath-Mistral-7B, where Full-Step-DPO shows a slight 0.8% drop in accuracy. However, this drop is notably smaller than 3.0% with DPO and the 3.7% with Step-DPO. These results demonstrate the superior stability and resilience of Full-Step-DPO, particularly in handling challenging mathematical reasoning tasks. More experimental results on additional datasets can be found in Appendix B.1.

Model	GK2023 (%)	OCW (%)
MetaMath-Mistral-7B	15.8	<b>10.7</b>
+ DPO	15.8 <sup>-0</sup>	7.7 <sup>-3.0</sup>
+ Step-DPO	15.1 <sup>-0.7</sup>	7.0 <sup>-3.7</sup>
+ Full-Step-DPO	<b>20.5</b> <sup>+4.7</sup>	<b>9.9</b> <sup>-0.8</sup>
Llama-3-8B-SFT	20.5	12.5
+ DPO	11.7 <sup>-8.8</sup>	9.9 <sup>-2.6</sup>
+ Step-DPO	19.7 <sup>-0.8</sup>	13.6 <sup>+1.1</sup>
+ Full-Step-DPO	<b>22.1</b> <sup>+1.6</sup>	<b>15.1</b> <sup>+2.6</sup>
DeepSeekMath-Base-SFT	30.4	19.1
+ DPO	31.2 <sup>+0.8</sup>	18.4 <sup>-0.7</sup>
+ Step-DPO	31.2 <sup>+0.8</sup>	18.0 <sup>-1.1</sup>
+ Full-Step-DPO	<b>31.7</b> <sup>+1.3</sup>	<b>20.2</b> <sup>+1.1</sup>
Qwen2-7B-SFT	33.0	15.8
+ DPO	8.8 <sup>-24.2</sup>	8.1 <sup>-7.7</sup>
+ Step-DPO	32.5 <sup>-0.5</sup>	15.8 <sup>-0</sup>
+ Full-Step-DPO	<b>33.5</b> <sup>+0.5</sup>	<b>18.4</b> <sup>+2.6</sup>

Table 2: Performance comparison on out-of-domain math problems.

### 4.4 Results on Various Verification Strategies

Figure 3 presents the performance of different verification strategies on GSM8K under two settings:  $K = 5$  and  $K = 15$ . We find that: (1) SC serves as a simple yet powerful validation method that significantly improves performance across all models. Even for the high-performing Qwen2-7B-SFT, which achieves an accuracy of 89.3% with greedy decoding, SC further improves the accuracy to 93% when the sampling size  $K = 15$ . This result is already comparable to GPT-4’s accuracy of 93.6%. (2) Compared to SC, BoN often achieves further improvements on MetaMath-Mistral-7B and Llama-3-8B-SFT. However, on the highly capable DeepSeekMath-Base-SFT and Qwen2-7B-SFT, BoN underperforms SC, indicating that the benefits of the reward model diminish for very strong baseline models. (3) SBS performs worse than both SC and BoN across most settings, yet consistently surpasses Greedy decoding, aligning with findings from previous studies (Yu et al., 2023a; Khalifa et al., 2023). This may be because, during the early stages of inference, the reward model struggles to effectively distinguish the correctness of steps.

### 4.5 Analysis of PRMs

As discussed in Section 3.3, the quality of the PRM may be influenced by the sampling number  $M$  and simulation number  $N$ . To assess this, we conducted a controlled experiment with fixed  $M = 32$  and varying  $N$ . The trained PRM is then used to guide the decoding of MetaMath-Mistral-7B-Full-Step-

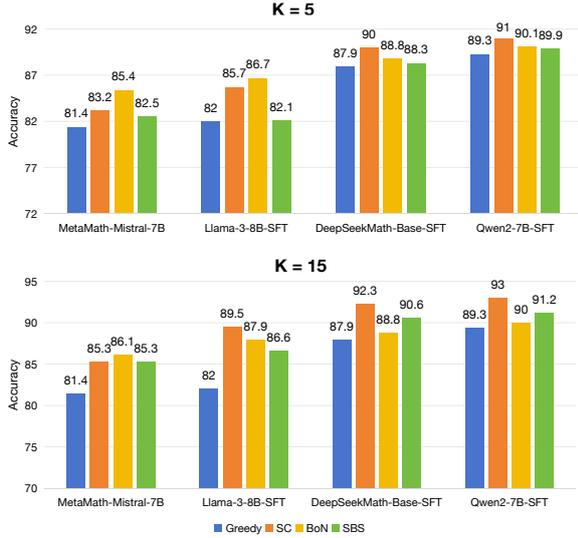


Figure 3: Performance comparison of various verifications on GSM8K, with all models trained using our Full-Step-DPO.

DPO using BoN strategy with  $K = 15$ .

As shown in Figure 4, the purple bars indicate the A100-Hours cost for constructing PRM training data, while blue and orange lines show GSM8K and MATH accuracy. When  $N = 1$ , which represents our proposed method, the model achieves competitive performance with 85.3% accuracy on GSM8K and 34.2% on MATH, while only requiring 8.5h with a time complexity of  $O(M)$ . As  $N$  increases, we observe that the performance fluctuation remains relatively minimal (within approximately 1% range), while the sampling cost grows substantially with a complexity of  $O(MNK)$ . This empirical evidence suggests that our approach with  $N = 1$  achieves a good balance between sampling efficiency and model performance. Additionally, we provide a comprehensive comparison between our PRM and other publicly available PRMs in Appendix B.2.

#### 4.6 Sensitivity of Hyperparameters

In step-wise DPO loss, the reward temperature  $\gamma$  reflects the level of trust in the PRM. As  $\gamma$  increases, the PRM model has a greater impact on the gradients. When  $\gamma \rightarrow 0$ , it indicates complete distrust in the PRM model, assigning equal weight coefficient to all steps, degrading in vanilla DPO. Conversely, when  $\gamma \rightarrow \infty$ , the loss function optimizes only the single step with the maximum or minimum reward in the solution, similar to Step-DPO. Figure 5 presents the accuracy of MetaMath-Mistral-7B-Full-Step-DPO

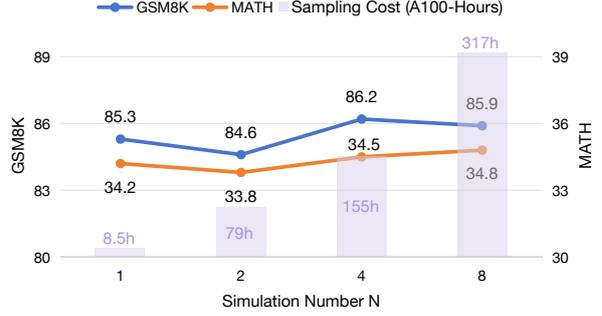


Figure 4: Accuracy of MetaMath-Mistral-7B-Full-Step-DPO using BoN decoding with  $K = 15$ . The PRM uses a fixed sampling number  $M = 32$ , while the simulation number  $N$  varies. Purple bars indicate the A100-Hours cost for constructing PRM training data.

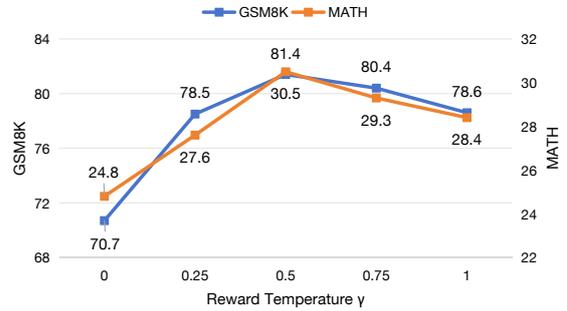


Figure 5: Accuracy of MetaMath-Mistral-7B-Full-Step-DPO with different reward temperature  $\gamma$ .

with different  $\gamma$  values. The experimental results indicate that introducing a PRM to weight the gradients indeed effectively enhances optimization efficiency and improves performance. Additionally, this experiment demonstrates that there is a sweet spot for the reward temperature  $\gamma$ ; excessively high or low  $\gamma$  will reduce accuracy.

## 5 Conclusion

In this work, we propose Full-Step-DPO, a novel framework for mathematical reasoning that optimizes each step in the entire reasoning chain using step-wise rewards. To achieve this, we train a self-supervised Process Reward Model to automatically score reasoning steps, eliminating reliance on external annotations. We also propose a novel Step-Wise DPO Loss that dynamically updates gradients based on the rewards for individual steps, enabling step-wise optimization and enhancing the reasoning ability of policy models. Experimental results on various benchmarks validate the effectiveness of Full-Step-DPO, paving the way for its application to other reasoning-intensive tasks.

## 611 Limitations

612 While we have conducted comprehensive exper-  
613 iments to demonstrate the effectiveness of Full-  
614 Step-DPO, several limitations remain. First, recent  
615 advancements suggest that generative reward mod-  
616 els outperform the discriminative reward model  
617 used in this work. Exploring how generative re-  
618 ward models can further enhance mathematical  
619 reasoning capabilities would be a valuable direc-  
620 tion for future research. Second, during preference  
621 data construction, the current strategy of selecting  
622 samples based on average reward is relatively sim-  
623 ple. Investigating more advanced sample selection  
624 strategies may lead to further improvements. Fi-  
625 nally, the step-wise DPO loss proposed in this paper  
626 is highly adaptable to other reasoning tasks, such  
627 as code generation. Conducting experiments on  
628 a broader range of tasks would provide additional  
629 evidence of the advantages of our method.

## 630 References

631 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
632 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
633 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
634 Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#).  
635 *arXiv preprint arXiv:2303.08774*.

636 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster,  
637 Marco Dos Santos, Stephen McAleer, Albert Q Jiang,  
638 Jia Deng, Stella Biderman, and Sean Welleck. 2023.  
639 [Llemma: An open language model for mathematics](#).  
640 *arXiv preprint arXiv:2310.10631*.

641 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
642 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
643 Huang, et al. 2023. [Qwen technical report](#). *arXiv*  
644 *preprint arXiv:2309.16609*.

645 Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai  
646 Fan. 2024. [Step-level value preference optimiza-](#)  
647 [tion for mathematical reasoning](#). *arXiv preprint*  
648 *arXiv:2406.10858*.

649 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
650 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
651 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
652 Nakano, et al. 2021a. [Training verifiers to solve math](#)  
653 [word problems](#). *arXiv preprint arXiv:2110.14168*.

654 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
655 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
656 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
657 Nakano, et al. 2021b. [Training verifiers to solve math](#)  
658 [word problems](#). *arXiv preprint arXiv:2110.14168*.

659 Rémi Coulom. 2006. [Efficient selectivity and backup](#)  
660 [operators in monte-carlo tree search](#). In *International*  
661 *conference on computers and games*, pages 72–83.  
662 Springer.

Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, 663  
Dan Jurafsky, and Douwe Kiela. 2024. [Kto: Model](#)  
664 [alignment as prospect theoretic optimization](#). *arXiv*  
665 *preprint arXiv:2402.01306*. 666

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul 667  
Arora, Steven Basart, Eric Tang, Dawn Song, and Ja- 668  
cob Steinhardt. 2021. [Measuring mathematical prob-](#) 669  
[lem solving with the math dataset](#). *arXiv preprint* 670  
*arXiv:2103.03874*. 671

Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren 672  
Etzioni, and Nate Kushman. 2014. [Learning to solve](#) 673  
[arithmetic word problems with verb categorization](#). 674  
In *Proceedings of the 2014 Conference on Empirical* 675  
*Methods in Natural Language Processing (EMNLP)*, 676  
pages 523–533. 677

Jie Huang, Xinyun Chen, Swaroop Mishra, 678  
Huaixiu Steven Zheng, Adams Wei Yu, Xiny- 679  
ing Song, and Denny Zhou. 2023. [Large language](#) 680  
[models cannot self-correct reasoning yet](#). *arXiv* 681  
*preprint arXiv:2310.01798*. 682

Fangkai Jiao, Chengwei Qin, Zhengyuan Liu, Nancy F 683  
Chen, and Shafiq Joty. 2024. [Learning planning-](#) 684  
[based reasoning by trajectories collection and](#) 685  
[process reward synthesizing](#). *arXiv preprint* 686  
*arXiv:2402.00658*. 687

Jean Kaddour, Joshua Harris, Maximilian Mozes, Her- 688  
bie Bradley, Roberta Raileanu, and Robert McHardy. 689  
2023. [Challenges and applications of large language](#) 690  
[models](#). *arXiv preprint arXiv:2307.10169*. 691

Muhammad Khalifa, Lajanugen Logeswaran, Moon- 692  
tae Lee, Honglak Lee, and Lu Wang. 2023. [Grace:](#) 693  
[Discriminator-guided chain-of-thought reasoning](#). In 694  
*Findings of the Association for Computational Lin-* 695  
*guistics: EMNLP 2023*, pages 15299–15328. 696

Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xi- 697  
anru Peng, and Jiaya Jia. 2024. [Step-dpo: Step-wise](#) 698  
[preference optimization for long-chain reasoning of](#) 699  
[llms](#). *arXiv preprint arXiv:2406.18629*. 700

Aitor Lewkowycz, Anders Andreassen, David Dohan, 701  
Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, 702  
Ambrose Slone, Cem Anil, Imanol Schlag, Theo 703  
Gutman-Solo, et al. 2022. [Solving quantitative rea-](#) 704  
[soning problems with language models](#). *Advances* 705  
*in Neural Information Processing Systems*, 35:3843– 706  
3857. 707

Minpeng Liao, Wei Luo, Chengxi Li, Jing Wu, and 708  
Kai Fan. 2024. [Mario: Math reasoning with code](#) 709  
[interpreter output—a reproducible pipeline](#). *arXiv* 710  
*preprint arXiv:2401.08190*. 711

Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri 712  
Edwards, Bowen Baker, Teddy Lee, Jan Leike, 713  
John Schulman, Ilya Sutskever, and Karl Cobbe. 714  
2023. [Let’s verify step by step](#). *arXiv preprint* 715  
*arXiv:2305.20050*. 716



827 Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, James Xu  
828 Zhao, Min-Yen Kan, Junxian He, and Michael Xie.  
829 2024b. [Self-evaluation guided beam search for rea-](#)  
830 [soning](#). *Advances in Neural Information Processing*  
831 *Systems*, 36.

832 Yifan Xu, Xiao Liu, Xinghan Liu, Zhenyu Hou, Yueyan  
833 Li, Xiaohan Zhang, Zihan Wang, Aohan Zeng,  
834 Zhengxiao Du, Wenyi Zhao, et al. 2024. [Chatglm-](#)  
835 [math: Improving math problem-solving in large lan-](#)  
836 [guage models with a self-critique pipeline](#). *arXiv*  
837 *preprint arXiv:2404.02893*.

838 Huaiyuan Ying, Shuo Zhang, Linyang Li, Zhejian Zhou,  
839 Yunfan Shao, Zhaoye Fei, Yichuan Ma, Jiawei Hong,  
840 Kuikun Liu, Ziyi Wang, et al. 2024. [Internlm-math:](#)  
841 [Open math large language models toward verifiable](#)  
842 [reasoning](#). *arXiv preprint arXiv:2402.06332*.

843 Fei Yu, Anningzhe Gao, and Benyou Wang. 2023a.  
844 [Outcome-supervised verifiers for planning in mathe-](#)  
845 [matical reasoning](#). *arXiv preprint arXiv:2311.09724*.

846 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu,  
847 Zhengying Liu, Yu Zhang, James T Kwok, Zhen-  
848 guo Li, Adrian Weller, and Weiyang Liu. 2023b.  
849 [Metamath: Bootstrap your own mathematical ques-](#)  
850 [tions for large language models](#). *arXiv preprint*  
851 *arXiv:2309.12284*.

852 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan  
853 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
854 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024.  
855 [Judging llm-as-a-judge with mt-bench and chatbot](#)  
856 [arena](#). *Advances in Neural Information Processing*  
857 *Systems*, 36.

## A Step-wise Beam Search

Unlike conventional beam search, which relies on token-level probabilities, our method integrates the reward model with an associated reranking criterion. This enables for step-wise beam search (SBS) (Yu et al., 2023a; Chen et al., 2024), effectively selecting the preferred solution path in mathematical reasoning, while incurring a lower computational cost compared to Monte Carlo Tree Search. Specifically, for each step  $t$ , suppose the sampling size is  $b_1$ , the policy model  $\pi_\theta$  produces a set of candidate steps  $\mathbb{S}^{(1:t+1)} = \{\mathcal{S}_i^{(1:t+1)}\}_{i=1}^{b_1}$ , where  $\mathcal{S}_i^{(1:t+1)} = [s_i^1, \dots, s_i^{t+1}]$  is the  $i$ -th partial solution up to step  $t + 1$ . Given the PRM  $\pi_r$  that can score each step, we select the top-scoring steps with beam size  $b_2$ . The algorithm is detailed in Algorithm 1. By focusing on the quality of each reasoning step rather than just the final solution, our method enhances the overall reasoning capabilities of the model.

---

### Algorithm 1 Step-wise Beam Search

---

```

1: Input: Math problem  $q$ , Sampling size  $b_1$ , Beam size  $b_2$ , Maximum step  $C$ 
2: Output: Best solution for  $q$ 
3: Models: Policy model  $\pi_\theta$  and PRM  $\pi_r$ 
4: function STEPLEVELBEAMSEARCH( $q, b_1, b_2, C$ )
5:   Initialize step sequences  $\mathbb{S} \leftarrow \{\}$ 
6:   Use  $\pi_\theta$  to sample initial steps  $\{s_1^1, \dots, s_{b_1}^1\}$ 
7:   Use  $\pi_r$  to score all initial steps  $\{r_1^1, \dots, r_{b_1}^1\}$ 
8:   Select top- $b_1$  steps and add to  $\mathbb{S}$ 
9:   Set current step counter  $t \leftarrow 1$ 
10:  while  $t < C$  do
11:    if All sequences in  $\mathbb{S}$  are complete then
12:      Break
13:    end if
14:     $\mathbb{S}_{\text{new}} \leftarrow \{\}$ 
15:     $\mathbb{R} \leftarrow \{\}$ ;
16:    for each solution  $\mathcal{S}^{(1:t)}$  in  $\mathbb{S}$  do
17:      for  $i = 1$  to  $b_1$  do
18:         $\mathcal{S}_i^{(1:t+1)} = \pi_\theta(\mathcal{S}^{(1:t)}; q)$ 
19:         $r_i^{(1:t+1)} = \pi_r(\mathcal{S}_i^{(1:t+1)}; q)$ 
20:         $\mathbb{S}_{\text{new}} \leftarrow \mathbb{S}_{\text{new}} + \{\mathcal{S}_i^{(1:t+1)}\}$ 
21:         $\mathbb{R} \leftarrow \mathbb{R} + \{r_i^{(1:t+1)}\}$ 
22:      end for
23:    end for
24:     $\mathbb{S}_{\text{new}} \leftarrow$  top- $b_2$  rewarded solutions in  $(\mathbb{S}_{\text{new}}, \mathbb{R})$ 
25:     $\mathbb{S} \leftarrow \mathbb{S}_{\text{new}}$ 
26:     $t \leftarrow t + 1$ ;
27:  end while
28:  return solution with highest final reward in  $\mathbb{S}$ 
29: end function

```

---

## B Additional Experiments

### B.1 More OOD Datasets

We evaluate our method on five additional OOD mathematical reasoning datasets as our testbed. As shown in Table 3, our Full-Step-DPO consistently improves performance across all datasets, demonstrating the effectiveness and generalization ability of our approach on OOD mathematical reasoning tasks.

- SVAMP (Patel et al., 2021) includes 1000 math questions of up to fourth grade difficulty. These questions can be solved by expressions requiring no more than two operators.

- AddSub (Hosseini et al., 2014) contains 395 math questions that involve addition and subtraction operations. 875 876
- ASDiv (Miao et al., 2021) contains 2215 English math questions of different problem types. Each question provides the corresponding equation and answer. 877 878
- GSM-IC2 and GSM-ICM (Shi et al., 2023) are mathematical reasoning datasets containing irrelevant conditions within the problem descriptions each consisting of 1000 problems. Problems in GSM-IC2 require two steps to solve, while problems in GSM-ICM require more than two steps to solve. 879 880 881

Model	SVAMP (%)	AddSub (%)	ASDiv (%)	GSM-IC2 (%)	GSM-ICM (%)	Average (%)
MetaMath-Mistral-7B	79.1	86.6	81.2	<b>77.9</b>	<u>76.5</u>	80.3
+ DPO	72.7	53.4	73.8	60.5	65.2	65.1
+ Step-DPO	<u>80.3</u>	86.9	<u>82.7</u>	76.4	76.3	<u>80.5</u>
+ Full-Step-DPO	<b>81.7</b>	<b>88.6</b>	<b>83.9</b>	75.9	<b>76.6</b>	<b>81.3</b>
Llama-3-8B-SFT	<u>82.8</u>	88.4	85.0	80.1	79.5	83.2
+ DPO	72.2	60.0	74.2	64.1	66.3	67.4
+ Step-DPO	81.7	<u>88.5</u>	<u>85.3</u>	<u>80.7</u>	<u>81.2</u>	<u>83.5</u>
+ Full-Step-DPO	<b>82.9</b>	<b>88.8</b>	<b>86.4</b>	<b>82.1</b>	<b>81.6</b>	<b>84.4</b>
DeepSeekMath-Base-SFT	84.2	<b>87.6</b>	<u>91.0</u>	85.4	85.2	86.7
+ DPO	85.1	86.6	90.6	85.2	85.0	86.5
+ Step-DPO	<u>85.3</u>	85.3	90.7	<u>85.9</u>	<b>86.2</b>	<u>86.7</u>
+ Full-Step-DPO	<b>85.9</b>	<u>86.6</u>	<b>91.2</b>	<b>87.8</b>	<u>85.3</u>	<b>87.4</b>
Qwen2-7B-SFT	88.7	<b>92.7</b>	91.6	<u>93.7</u>	91.6	<u>91.7</u>
+ DPO	23.5	25.6	27.0	30.1	29.5	27.1
+ Step-DPO	88.1	92.2	<u>91.8</u>	<b>93.8</b>	91.9	91.6
+ Full-Step-DPO	<b>89.5</b>	<b>93.1</b>	<b>92.4</b>	93.6	<b>92.7</b>	<b>92.3</b>

Table 3: Performance comparison on five additional OOD math problems. The best and second-best results are highlighted in bold and underlined, respectively.

## B.2 Performance Comparison of Different PRMs 882

We compare our PRM with Math-Shepherd-PRM-7B using MetaMath-Mistral-7B as the base model under the Best-of-8 decoding strategy. During decoding, we sampled eight responses with temperature = 0.7. For both PRMs, each response was scored by the minimum reward among its reasoning steps, and the response with the highest score was selected as the final output. 883 884 885 886

As shown in Table 4, our PRM achieves 84.3% accuracy on GSM8K and 33.2% on MATH, outperforming Math-Shepherd-PRM-7B (82.4% on GSM8K, 32.6% on MATH). While these results suggest the effectiveness of our approach, we acknowledge that the comparison may not be entirely comprehensive due to differences in backbone architecture and training data between the PRMs. Nevertheless, these findings, combined with our analysis in Section 4.5, demonstrate that our PRM can achieve competitive performance while significantly reducing the computational cost of PRM construction. 887 888 889 890 891 892

Model	MATH	GSM8K
MetaMath-Mistral-7B	28.2	77.7
+ Self-Consistency	32.4	81.3
+ Math-Shepherd-PRM-7B	32.6	82.4
+ Our-PRM-7B	<b>33.2</b>	<b>84.3</b>

Table 4: Performance comparison of different scoring strategies using Best-of-8 decoding on MATH and GSM8K benchmarks.

893  
894  
895  
896  
897  
898  
899

### B.3 Case Study

We conducted a detailed case study using a sample problem from the AddSub dataset, as shown in Table 5. We compared outputs from DeepSeekMath-Base-SFT model and its Step-DPO and Full-Step-DPO variants. While both the base SFT model and Step-DPO model made incorrect reasoning by including the cost of unbought CDs, our Full-Step-DPO model generated correct reasoning steps. Our PRM scores accurately reflect the correctness of each reasoning step, with higher scores for correct steps and lower scores for incorrect ones.

<b>Question</b>	For his car, Mike spent \$118.54 on speakers and \$106.33 on new tires. Mike wanted 3 CD's for \$4.58 but decided not to. In total, how much did Mike spend on car parts?
<b>Gold Answer</b>	224.87
<b>SFT Response</b>	<b>Step 1:</b> Mike spent 118.54 on speakers and 106.33 on new tires. So the total cost of these two items is $118.54 + 106.33 = 224.87$ . <b>Score: 0.9348</b>
	<b>Step 2:</b> Mike wanted to buy 3 CD's for \$4.58 each, but he decided not to. So the total cost of the CD's he didn't buy is $3 * 4.58 = 13.74$ . <b>Score: 0.3646</b>
	<b>Step 3:</b> To find out how much Mike spent in total on car parts, we need to add the total cost of the speakers and tires to the total cost of the CD's he didn't buy. So the total cost is $224.87 + 13.74 = 238.61$ . The answer is: 238.61 <b>Score: 0.0427</b>
<b>Step-DPO Response</b>	<b>Step 1:</b> First, Mike spent \$118.54 on speakers and \$106.33 on new tires. So the total cost of these two items is $118.54 + 106.33 = 224.87$ . <b>Score: 0.9406</b>
	<b>Step 2:</b> Mike also wanted 3 CD's, but he decided not to buy them. Each CD costs \$4.58, so the total cost of the CD's he didn't buy is $3 * 4.58 = 13.74$ . <b>Score: 0.3036</b>
	<b>Step 3:</b> To find out how much Mike spent in total on car parts, we need to add the total cost of the speakers and tires to the total cost of the CD's he didn't buy. So, the total cost is $224.87 + 13.74 = 238.61$ . The answer is: 238.61 <b>Score: 0.0445</b>
<b>Full-Step-DPO Response</b>	<b>Step 1:</b> Mike spent \$118.54 on speakers and \$106.33 on new tires. So the total amount he spent on these two items is $118.54 + 106.33 = 224.87$ . <b>Score: 0.9313</b>
	<b>Step 2:</b> Mike wanted to buy 3 CD's for \$4.58 each, but he decided not to. So the total amount he would have spent on the CD's is $3 * 4.58 = 13.74$ . <b>Score: 0.7408</b>
	<b>Step 3:</b> To find out how much Mike spent in total on car parts, we need to add the total amount he spent on speakers and tires to the amount he would have spent on CD's. So the total amount is $224.87 + 0 = 224.87$ . The answer is: 224.87 <b>Score: 0.7363</b>

Table 5: Comparison of reasoning steps and their corresponding PRM scores across different models for AddSub dataset. The base model is DeepSeekMath-Base-SFT.