STEP-CONTROLLED DPO: LEVERAGING STEPWISE ERRORS FOR ENHANCING MATHEMATICAL REASON ING OF LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Direct Preference Optimization (DPO) has proven effective at improving the performance of large language models (LLMs) on downstream tasks such as reasoning and alignment. In this work, we propose Step-Controlled DPO (SCDPO), a method for automatically providing stepwise error supervision by creating negative samples of mathematical reasoning rationales that start making errors at a specified step. By applying these samples in DPO training, SCDPO can better align the model to avoid reasoning errors and output accurate reasoning steps. Qualitative analysis of the credit assignment of SCDPO and DPO demonstrates the effectiveness of SCDPO at identifying errors in mathematical solutions. We then apply SCDPO to an InternLM2-20B model, resulting in a 20B model that achieves competitive scores of 88.5% on GSM8K and 58.1% on MATH, rivaling all other open-source LLMs, showing the great potential of our method. The code, models and data are released to inspire future work.

- 1 Ιντροριά
- 027 028

025 026

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

Recently, Direct Preference Optimization (DPO; Rafailov et al. (2024b)) has emerged as a popular choice for aligning large language models (LLMs) with relative feedback to improve the quality of generated text. Prior works Christiano et al. (2023); Pal et al. (2024); Xu et al. (2024) have demonstrated that reinforcement learning algorithms and DPO can improve the mathematical reasoning abilities of LLMs, making the generated reasoning process more controllable. The final answer to a mathematical problem serves as a natural way to judge the quality of the model's response, since a mathematical problem typically has a single correct answer. As a result, the responses producing the correct final answers are desirable and can serve as the preferred samples, while the ones reaching incorrect final answers are undesirable and can serve as the dispreferred samples.

However, solutions to a mathematical problem can be diverse, with many different reasoning paths arriving at the correct final answer and many subtle ways to make mistakes. Determining the pre-ferred and dispreferred responses based on the final answer is coarse and may be inadequate for capturing *the intricacies of the multi-step mathematical reasoning process*. Previous studies intro-duce process supervision Lightman et al. (2023), but it requires large amounts of meticulous and expensive human annotation and only applies to traditional RL algorithms.

In this paper, we show how to automatically provide explicit stepwise preference supervision by 044 generating diverse dispreferred solutions that start making errors at a specific step. We propose *Step*-045 Controlled DPO (SCDPO), an simple yet effective algorithm that introduces stepwise supervision 046 without necessitating extra human annotation. This approach starts with a model finetuned with 047 question-solution pairs and possessing initial math-solving capabilities, which is used to generate 048 solutions to a set of math problems. We choose the solutions whose final answers match those of the ground truth. We take each of these correct solutions and start generating with the model via modulating the hyperparameter of the model, i.e., increasing the temperature of the final softmax 051 function, from various intermediate steps of that solution, and retain the samples where the final answer is incorrect. In this way, the steps before the intermediate step are the same as the original 052 correct solution, while the steps after are the ones with possible errors. During DPO training, the correct solutions are the preferred samples, and they are paired with the wrong solutions generated



Figure 1: Demonstration and example of the step-controlled data generation process. **a.** Stepcontrolled data generation. First, a solution reaching the correct final answers is collected, which we denote as $a_i^{(\text{pre})}$. Then, erroneous solutions that reach incorrect final answers are generated, starting from intermediate steps of $a_i^{(\text{pre})}$, creating dispreferred solutions $a_{i1}^{(\text{dis-sc})}$, $a_{i2}^{(\text{dis-sc})}$, and $a_{i3}^{(\text{dis-sc})}$. These dispreferred solutions share the steps before the intermediate steps with $a_i^{(\text{pre})}$. The temperature of the newly generated steps gradually increases with each step to make the generation more erroneous. **b.** An example of a pair of preferred and dispreferred solutions. The dispreferred solution starts making errors after a particular intermediate step.

087

090

091

092

094

096

098

099

100

077

in this way, with the question and the steps before the intermediate step as the prompts. These stepcontrolled training samples help models learn detailed reasoning abilities and are mixed with naive DPO training data produced by only checking the final answer, which optimizes the general form of the solution.

Our contributions are as follows:

- We introduce SCDPO, a method that automatically provides explicit stepwise supervision to enhance mathematical abilities of LLMs.
- We conduct pilot experiments on chain-of-thought and code-integrated solutions, showing that SCDPO can effectively improve mathematical problem-solving performance of three different SFT models. We also conduct qualitative analysis of credit assignment of SCDPO.
- Using SCDPO, we finetune an InternLM2-20B model, which reaches 88.5% on GSM8K Cobbe et al. (2021) and 58.1% on MATH Hendrycks et al. (2021), demonstrating the great potential of our method.
- 102 103

2 STEP-CONTROLLED DPO PIPELINE

104 105

In this section, we introduce Step-Controlled DPO (SCDPO), a pipeline for automatically generating
 preferred and dispreferred responses to math problems, with annotations of erroneous solving steps,
 and using these responses in DPO training to enhance the mathematical reasoning abilities of LLMs.

Our method consists of two stages: step-controlled data generation, and step-aware DPO training. The two stages construct a feedback-alignment framework that is both effective and cost-efficient.

Initial Model. Our method starts with an initial model, denoted as π_{SFT} , which has been finetuned with question-solution pairs from math datasets such as GSM8K and MATH. When prompted with a math problem q, π_{SFT} is able to generate a step-by-step solution, denoted as a. a can be broken down into a sequence of reasoning steps, for example, $a = (t_0, \ldots, t_m)$. Here, t_i $(i = 0, \ldots, m)$ represents either a code reasoning step or a natural language reasoning step within a. For Chainof-Thought solutions, the reasoning steps are separated by "\n". In code-integrated solutions, the reasoning steps are separated by special tokens as described in Wang et al. (2023a).

117 118 119

120

2.1 STEP-CONTROLLED DATA GENERATION

The data we collect is in two parts: naive DPO data D_{naive} and Step-Controlled DPO data D_{SC} .

122 Generation of D_{naive} . D_{naive} contains pairs of preferred-dispreferred samples, used to optimize the general form of the solution. To create D_{naive} , we prompt π_{SFT} with math questions in the training 123 sets of GSM8K and MATH. Each question is presented to π_{SFT} multiple times and various solutions 124 are generated, with a temperature of 1. If a solution reaches the same final answer as the ground 125 truth, and no errors or adjustments occur at any of the reasoning steps (we detect these by looking 126 for strings like "error" or "apologies"), the solution is seen as preferred, while the solutions that 127 reach answers different from the ground truth are considered dispreferred. To find out the frequency 128 of incorrect solution process reaching the correct final answer, we randomly sampled 87 solutions 129 that reach correct final answers, and found that of the 369 reasoning steps in these solutions, only 130 2 contain errors, which is a very small percentage (about 0.5%). This demonstrates that, in most 131 cases, a correct final answer indicates correct intermediate steps. The solution generation of each 132 question stops when at least one preferred solution and one dispreferred solution are generated, or 133 the number of solutions generated reaches an upper limit of 100. We use questions from the training sets of the GSM8K and MATH datasets for solution sampling, and repeated sampling ensures that 134 99.8% of the questions in the GSM8K training set and 91.8% of the questions in the MATH training 135 set yield at least one positive sample. The resulting data can be expressed as: 136

- 137
- 138
- 139 140

 $D_{\text{naive}} = \{ (q_i, a_i^{(\text{pre})}, a_i^{(\text{dis})}) : i = 1, \dots, N_{\text{naive}} \}$ (1)

Here, q_i denotes the *i*th question, while $a_i^{(\text{pre})}$ and $a_i^{(\text{dis})}$ represent the preferred and dispreferred solution to the *i*th question.

Generation of D_{SC} . In order to generate solutions with stepwise error information for DPO train-144 ing, we propose a method to automatically generate training data with errors starting to occur at 145 a controlled step. The process is demonstrated in Fig. 1. We first take a preferred solution from 146 D_{naive} , denoted as $a_i^{(\text{pre})} = (t_0^{(\text{pre})}, \dots, t_k^{(\text{pre})}, t_{k+1}^{(\text{pre})}, \dots, t_{m_i}^{(\text{pre})})$. Here, $t_k^{(\text{pre})}$ is a random intermediate step within $a_i^{(\text{pre})}$. As $a_i^{(\text{pre})}$ is a correct solution, $t_0^{(\text{pre})}, \dots, t_{m_i}^{(\text{pre})}$ can all be seen as correct steps. As shown in Fig. 1 a, to create a solution with errors occurring after step k, we present π_{SFT} with 147 148 149 sequence $(q_i, t_0^{(\text{pre})}, \ldots, t_k^{(\text{pre})})$, and raise the temperature of the final softmax function to affect the 150 generation quality, increasing the occurrence of errors in the following steps. Raising the tem-151 perature causes the model performance to become unstable and erroneous. We observe that when 152 the temperature is instantly raised and remains at a high value, the model can generate garbled 153 strings as errors accumulate, which does not represent any reasoning mistakes and contains no valu-154 able information. To avoid this, we adopt a gradually increasing temperature, which initially starts at 1.1, and increases by 0.05 with each generated step, until the generation ends or the tempera-156 ture reaches 1.4. This setting empirically reduces the frequency of the occurrence of garbled text, 157 while increasing the error rate and diversity of generated errors. We generate the steps following 158 $(q_i, t_0^{(\text{pre})}, \dots, t_k^{(\text{pre})})$ multiple times, until one reaching an incorrect answer is generated. Appending 159 the generated steps to $(t_0^{(\text{pre})}, \ldots, t_k^{(\text{pre})})$, we get a dispreferred solution with step-controlled error, denoted as $a_{ik}^{(\text{dis-sc})} = (t_0^{(\text{pre})}, \ldots, t_k^{(\text{pre})}, t_{k+1}^{(\text{dis-sc})}, \ldots, t_{n_i}^{(\text{dis-sc})})$, where the sequence $(t_{k+1}^{(\text{dis-sc})}, \ldots, t_{n_i}^{(\text{dis-sc})})$ is erroneous. An example is presented in Fig. 1 b. The resulting data can be expressed as: 160 161

 $D_{\text{SC}} = \{(q_i, a_i^{(\text{pre})}, a_{ik}^{(\text{dis-sc})}) : i = 1, \dots, N_{\text{SC}}\}$ (2)

Here, q_i denotes the *i*th question, while $a_i^{(\text{pre})}$ is the preferred solution, and $a_{ik}^{(\text{dis-sc})}$ is the dispreferred solution with step-controlled error that occurs after $t_k^{(\text{pre})}$. N_{SC} is the number of questions in D_{SC} , while m_i is the index of the last step of $a_i^{(\text{pre})}$.

2.2 STEP-CONTROLLED DPO TRAINING

Having collected D_{naive} and D_{SC} , we apply them to DPO training. D_{naive} serves to regulate the general form of solutions, while D_{SC} supervises the model's reasoning on a step level. During DPO training, samples in D_{naive} and D_{SC} are mixed together randomly, and the DPO loss is applied to each sample. For samples from D_{naive} , the loss is applied to all steps in the preferred and dispreferred solutions, which can be written as:

$$\mathcal{L}_{\text{naive}}(\pi_{\theta}; \pi_{\text{SFT}})$$

$$= -\mathbb{E}_{(q_i, a_i^{(\text{pre})}, a_i^{(\text{dis})}) \sim \mathcal{D}_{\text{naive}}} [\log \sigma(\beta \log \frac{\pi_{\theta}(a_i^{(\text{pre})} | q_i)}{\pi_{\text{SFT}}(a_i^{(\text{pre})} | q_i)}$$

$$-\beta \log \frac{\pi_{\theta}(a_i^{(\text{dis})} | q_i)}{\pi_{\text{SFT}}(a_i^{(\text{dis})} | q_i)})]$$
(3)

For a pair of preferred and dispreferred solutions in D_{SC} where the erroneous steps are generated starting from the kth step of the preferred solution, the preferred solution can be denoted as $a_i^{(pre)}$, and the dispreferred solution can be denoted as $a_{ik}^{\text{(dis-sc)}}$. The first k reasoning steps are shared between the pair of solutions. The erroneous steps after the kth step in $a_{ik}^{\text{(dis-sc)}}$ is denoted as $(t_{k+1}^{\text{(dis-sc)}}, \ldots, t_{n_i}^{\text{(dis-sc)}})$, while the correct steps in $a_i^{(\text{pre})}$ after the kth step is denoted as $(t_{k+1}^{(\text{pre})}, \ldots, t_{m_i}^{(\text{pre})})$. SCDPO directly contrast between the steps in $a_{ik}^{(\text{pre})}$ and $a_{ik}^{(\text{dis-sc})}$ after the kth step, applying the DPO loss only on the different steps.

 $\mathcal{L}_{SC}(\pi_{\theta};\pi_{SFT}) =$

$$-\mathbb{E}_{(q_i, a_i^{(\text{pre})}, a_{ik}^{(\text{dis-sc})}) \sim \mathcal{D}_{\text{SC}}}[\log \sigma((\sum_{j=k+1}^{m_i} \beta \log \frac{\pi_{\theta}(t_j^{(\text{pre})} | q_i, t_{< j})}{\pi_{\text{SFT}}(t_j^{(\text{pre})} | q_i, t_{< j})}) - (\sum_{j=k+1}^{n_i} \beta \log \frac{\pi_{\theta}(t_j^{(\text{dis-sc})} | q_i, t_{< j})}{\pi_{\text{SFT}}(t_j^{(\text{dis-sc})} | q_i, t_{< j})}))]$$

Combining \mathcal{L}_{naive} and \mathcal{L}_{SC} , the final loss function of Step-Controlled DPO is as follows:

$$\mathcal{L}_{\rm SCDPO} = \mathcal{L}_{\rm naive} + \mathcal{L}_{\rm SC} \tag{5}$$

(4)

In this way, \mathcal{L}_{naive} optimizes the general form of the solution, while \mathcal{L}_{SC} focuses on detailed reasoning steps, thus improving the model's accuracy in solving mathematical problems.

THEORETICAL EXPLANATION OF STEP-CONTROLLED DPO

Theoretical Insight. In this section, we provide some theoretical insights into why SCDPO can effectively enhance the reasoning ability of LLMs. As explained in Rafailov et al. (2024a), the DPO loss can be cast into token-level MDP (Markov Decision Process; Puterman (1994)). Similarly,

216	To find out how many fireflies remained, we need to follow these	To find out how many fireflies remained, we need to follow these
217	steps:	steps:
218	1 Start with the initial count of fireflies: 3.	1. Start with the initial count of fireflies: 3.
219	2 <mark>. Add 4 less</mark> than a dozen more fireflies: \(\text{dozen} = 12 \),	2. Add 4 less than a dozen more fireflies: \(\text{dozen} = 12 \),
220	so \(4 \times \text{(less than a dozen)} = 4 \times (12 - 4) \).	so \(4 \times \text{(less than a dozen)} = 4 \times (12 - 4) \).
221	3. Find the total number of fireflies before two flew away.	3. Find the total number of fireflies before two flew away.
222	4. Subtract 2 because two of the fireflies flew away.	4. Subtract 2 because two of the fireflies flew away.

Figure 2: Credit assignment of part of a solution for a GSM8K problem. Each token is colored corresponding to the DPO implicit reward as expressed in Eq. 6 (darker is higher). The left is the credit assignment of SCDPO, which correctly highlights the error -4 less than a dozen is not 4 times (12 - 4), while the credit assignment of DPO on the right fails to highlight it. 227



Figure 3: Credit assignment of part of a solution for a MATH problem. Each token is colored corresponding to the DPO implicit reward as expressed in Eq. 6 (darker is higher). The left is the credit assignment of SCDPO, which correctly highlights the error – as the original question was "Find the remainder when $8 \cdot 10^{18} + 1^{18}$ is divided by 9", the remainders of the terms 8, 10^{18} , and 1^{18} should not be summed, while the credit assignment of DPO on the right fails to highlight the error.

we can also interpret DPO as a step-level MDP. As presented in Eq. 4, $\beta \log \frac{\pi_{\theta}(t_j^{\text{(pre)}}|q_i, t_{<j})}{\pi_{\text{SFT}}(t_i^{\text{(pre)}}|q_i, t_{<j})}$ and

 $\beta \log \frac{\pi_{\theta}(t_j^{(\text{dis-sc})}|q_i, t_{< j})}{\pi_{\text{SFT}}(t_j^{(\text{dis-sc})}|q_i, t_{< j})}$ represent the reward of a single preferred or dispreferred step. For naive DPO, 246 247 all steps in the preferred and dispreferred solutions have their rewards affecting the loss. However, 248 many steps in the dispreferred solution are actually correct, as the error often occurs in a later step. 249 Step-Controlled DPO reduces the range of steps, starting from the (k + 1)th step, from which the 250 dispreferred steps are more likely to be erroneous due to the raised sampling temperature. The focus 251 of the optimization is thus cast on the errored steps rather than the whole solution, letting the model 252 learn more detailed reasoning abilities.

253 **Oualitative Evaluation of Credit Assignment of SCDPO.** We perform qualitative evaluation of 254 credit assignment on two models trained with SCDPO and DPO respectively. For a sequence of 255 tokens $\mathbf{x} = (x_0, \dots, x_m)$, where x_i is the *i*th token in the sequence, we denote all the tokens before 256 x_i as s_i , written as $s_i = (x_0, \ldots, x_{i-1})$. As introduced in recent research Rafailov et al. (2024a), 257 the DPO implicit reward can be expressed as follows:

224

225

226

237

238

239

240

241

242 243

244 245

$$r(\mathbf{s}_i, x_i) = \beta \log \pi(x_i | \mathbf{s}_i) - \beta \log \pi_{\text{SFT}}(x_i | \mathbf{s}_i)$$
(6)

261 Here $r(\mathbf{s}_i, x_i)$ denotes the DPO implicit reward of token x_i , which is the value we visualize as the 262 background color of the token. A darker color represents a higher reward value. As demonstrated in Fig. 2 and Fig. 3, when presented with an incorrect reasoning step, SCDPO more accurately 264 identifies the incorrect tokens compared to DPO. Fig. 2 shows part of a solution for a GSM8K 265 question. In step 2, the solution incorrectly interprets "4 less than a dozen" as " $4 \times (12 - 4)$ ", when 266 it should have been "(12 - 4)". The SCDPO model correctly highlights " $4 \times (12 - 4)$ ", while the DPO does not. Fig. 3 shows part of a solution for a MATH question. The solution sums the terms 267 in the expression when two of the terms should have been multiplied. SCDPO correctly highlights 268 the incorrect solution, while DPO does not. These examples show that the stepwise supervision 269 provided in SCDPO results in a better token-level understanding of reasoning errors.

Model	Size			Ε	nglish				C	hinese	
Houti	5124	GSM8	K MATI	носм	/ hung- arian	Mathe- matics	-SVA- MP	Simu eq	l-APE- 210K	CMA- TH	- MGSM zh
			Closed-	Source	Models						
GPT-3.5	-	80.8	34.1	-	41	-	-	-	-	73.8	-
GPT-4	-	93.6	<u>53.6</u>	30.1	92	-	-	-	84.2	89.3	-
GPT-4 Code Interpreter	-	97.0	69.7	-	-	-	-	-	-	-	-
GLM-4 ¹	-	91.8	49.0	-	<u>75</u>	-	-	-	93.5	<u>89.0</u>	-
			Open-S	ource N	Models						
Qwen2	7B	85.7	52.9	10.7	56	51.4	86.3	83.6	54.2	73.8	58.0
Math-Shepherd	7B	84.1	33.0	12.5	46	36.6	81.8	84.6	45.9	68.8	67.6
DeepSeekMath-RL	7B	86.7	58.8	22.1	55	57.4	86.7	69.6	71.9	87.6	78.4
SVPO	7B	81.7	59.5	34.2	-	-	-	-	-	-	-
InternLM2-Math	20B	80.7	54.3	12.9	66	41.1	83.4	55.6	64.3	69.0	58.4
MathGenie	20B	87.7	55.7	23.5	69	85.1	87.3	88.5	-	-	-
ChatGLM3-32B	32B	82.6	40.6	-	73	-	-	-	89.4	85.6	-
ToRA	34B	80.7	50.8	5.5	-	77.9	80.5	50.2	-	53.4	41.2
MAmmoTH	70B	76.9	41.8	-	-	65.4	84.3	51.8	-	-	-
MathCoder	70B	83.9	45.1	-	-	74.4	84.9	77.0	-	-	-
InternLM2-SFT	20B	86.4	55.8	21.6	71	84.0	86.9	91.2	77.1	88.4	74.8
InternLM2-SFT-DPO	20B	87.0	57.6	25.5	<u>74</u>	85.6	<u>89.7</u>	<u>92.6</u>	78.7	<u>89.9</u>	76.0
InternLM2-SFT-	20B	88.2	57.5	24.5	73	86.3	88.9	91.1	78.8	89.3	76.0
DPO(d-e)											
InternLM2-SFT-SCDPO	20B	88.5	58.1	<u>29.4</u>	78	87.5	90.2	93.6	<u>79.3</u>	90.3	80.4

Table 1: Performance of open-source and closed-source models on seven English datasets, GSM8K, MATH, OCW, hungarian, Mathematics, SVAMP and Simuleq, and three Chinese datasets, APE210K, CMATH, and MGSM-zh. All results reported are based on greedy decoding. The best models are marked in **bold**, and the second best models are <u>underlined</u>. Our 20B model trained on SCDPO outperforms SFT and naive DPO on all 10 datasets, demonstrating performance rivaling all other open-source models of similar scales.

302 303 304

305

4 EXPERIMENTS

In this section, we first train a 20B model using SCDPO, reaching a performance rivaling all other models of similar scale. Then, we perform a comprehensive empirical comparison between SCDPO and DPO on three kinds of Mistral-7B SFT models. We also present ablation studies to further explain the design of increasing temperature during the generation of erroneous steps and combining D_{naive} with D_{SC} during training.

311

312 4.1 20B MODEL TRAINED WITH SCDPO 313

Training Data. We collect solutions for questions in the training set of APE210K Zhao et al. (2020),
GSM8K and MATH from GPT-4 Code Interpreter. Combining 169K samples from APE210K, 34K
from GSM8K and 47K from MATH, we get an SFT dataset of 250K question-solution pairs. The
SCDPO and DPO training data is collected as described before in Sec. 2.1. During sampling, top-p
is set to 1 and top-k is set to -1 to consider all tokens. The training data for SCDPO contains 13K
samples from GSM8K, 46K samples from MATH, and 29K samples from APE210K.

Training Settings. We use InternLM2-20B Cai et al. (2024) as the foundation model, as it has demonstrated high performance in previous works Lu et al. (2024); Cai et al. (2024), even surpassing larger models such as Mixtral-8x7B Jiang et al. (2024) and Llama2-70B Touvron et al. (2023) in some cases. In the SFT stage, we finetune the model with a learning rate of 1.0×10^{-5} for 3 epochs, with a context length of 2048 tokens. In DPO and SCDPO training, we use a learning rate

Method	GSM8K	MATH	OCW	hungarian	Mathematics	SVAMP	Simuleq
			Mistral-71	B-Ours			
SFT (Baseline)	76.8	43.2	21.7	52	69.8	81.3	73.9
SFT-continued	76.3	43.9	18.8	55	70.3	80.8	74.5
SFT+DPO	78.8	45.1	18.4	56	74.8	81.0	74.9
SFT+DPO(d-e)	79.0	45.7	18.0	59	74.4	79.2	73.2
SFT+DPO+SC	80.1	47.7	22.4	61	76.5	82.3	79.0
		Ν	letaMath-N	listral-7B			
SFT (Baseline)	77.7	28.2	12.5	33	33.9	80.0	68.5
SFT-continued	76.8	28.5	13.2	35	33.6	80.3	69.1
SFT+DPO	81.0	28.7	14.0	34	33.8	81.0	68.3
SFT+DPO(d-e)	81.4	29.0	14.7	38	34.3	80.9	70.6
SFT+DPO+SC	81.7	29.3	15.4	42	35.0	81.6	73.2
		М	athCoder-N	Aistral-7B			
SFT (Baseline)	78.1	39.3	12.9	62	70.4	79.4	80.5
SFT-continued	78.2	40.3	12.5	65	71.2	77.3	80.7
SFT+DPO	79.2	42.9	14.3	65	74.9	85.4	81.3
SFT+DPO(d-e)	78.3	41.1	14.7	68	74.9	84.9	82.3
SFT+DPO+SC	80.4	43.4	15.7	70	75.4	85.4	83.1

Table 2: Effect of using Step-Controlled DPO (SCDPO) on three different SFT models: Mistral-7B-Ours, MetaMath-Mistral-7B and MathCoder-Mistral-7B. "(d-e)" denote the DPO baseline using the same amount of data as SCDPO. In all three cases, SCDPO outperforms the starting SFT model, continue pretraining on correct samples, naive DPO, and naive DPO with equal amount of data.

Model	GSM8K	MATH
Mistral-7B-Ours-SFT	76.8	43.2
Mistral-7B-Ours-SCDPO (temperature=1.0)	78.6	45.9
Mistral-7B-Ours-SCDPO (temperature=1.3)	80.0	45.9
Mistral-7B-Ours-SCDPO (ascending temperature)	80.1	47.7

Table 3: Pilot experiments of using different temperatures when generating error steps. When temperature equals 1.0, the errors are not diverse enough. When temperature equals 1.3, the model generates unintelligible strings due to accumulated errors. The design of ascending temperature offers more diversity while avoids generating meaningless errors, resulting in the best performance.

361

345

346

347

348

of 1.5×10^{-7} to train the SFT model for 2 epochs, with a context length of 1024 and β set to 0.1. The models are trained on 16 NVIDIA A800 80GB GPUs with a batch size of 64.

364 Evaluation Datasets. Ten representative mathematical datasets are used in evaluating the models: 365 GSM8K Cobbe et al. (2021), MATH Hendrycks et al. (2021), OCWCourses (OCW) Lewkowycz 366 et al. (2022), Hungarian National Exams (hungarian) Paster (2023), Mathematics Saxton et al. 367 (2019), SVAMP Patel et al. (2021), Simuleq Kushman et al. (2014), APE210K Zhao et al. (2020), 368 CMATH Wei et al. (2023b) and MGSM-zh Shi et al. (2023). The first seven datasets consist of En-369 glish math questions, while the last three consist of Chinese math questions. The evaluation datasets 370 contain a wide range of problem types, covering mathematical problems from grade-school level to college level, comprehensively evaluating the models' mathematical reasoning abilities. We use 371 greedy decoding for all evaluations. 372

Baselines. We compare our 20B models with powerful closed-source models such as GPT-374
3.5 (Brown et al., 2020), GPT-4 OpenAI et al. (2024), GPT-4 Code Interpreter OpenAI et al. (2024) and GLM-4², as well as open-source models such as MARIO Liao et al. (2024), Qwen2 Yang et al. (2024), Math-Shepherd Wang et al. (2024), SeaLLM-v2 Nguyen et al. (2024), DeepSeekMath-

³⁷⁷

²https://open.bigmodel.cn/dev/api#glm-4

Figure 4: Accuracy of Mistral-7B-Ours (SFT) on GSM8K and MATH when temperature is set at different values.

400 401

378

379380381382

384

385 386

391

392 393

397

398 399

RL Shao et al. (2024), SVPO Chen et al. (2024), Skywork-13B-Math Yang et al. (2023a), InternLM2-Math³ Ying et al. (2024), MathGenie Lu et al. (2024), ChatGLM3-32B-RFT-DPO Xu et al. (2024), Yi-Chat Yi (2023), ToRA Gou et al. (2024), MAmmoTH Yue et al. (2023), Math-Coer Wang et al. (2023a) and WizardMath Luo et al. (2023).

406 Main Results. Tab. 1 displays our main results, as well as various closed-source and open-source 407 baselines. Our model achieves a score of 88.5% on GSM8K, 78 on hungarian, 87.5% on Mathemat-408 ics, 90.2% on SVAMP, 93.6% on Simuleq, 90.3% on CMATH, and 80.4% on MGSM-zh, surpassing 409 all models with published parameters, and obtaining second-best scores among open-source models 410 on APE210K. Our model obtains a score of 58.1% on MATH, which is close to the best and second-411 best open-source score of 59.5% and 58.8%. While our model rivals the performance of GPT-3.5 412 on GSM8K and MATH, and surpasses GPT-4 and GLM-4 on MATH, it still underperforms GPT-4 413 Code Interpreter on GSM8K and MATH, and GLM-4 on APE210K.

414 Compared to InternLM2-SFT, InternLM2-SFT-SCDPO consistently increases the score on each of 415 the five datasets by approximately 2% to 3%. Compared to both InternLM2-SFT-DPO, which uses 416 the D_{naive} part of InternLM2-SFT-SCDPO's training data, and InternLM2-SFT-DPO_(data-equal), which 417 uses about the same amount of training data as InternLM2-SFT-SCDPO, InternLM2-SFT-SCDPO 418 consistently achieves the best performance across all five datasets, highlighting the effectiveness of 419 SCDPO in enhancing mathematical problem-solving abilities.

420 421

422

430

431

4.2 COMPARISON USING DIFFERENT STARTING 7B MODELS

We validate the generalizability of SCDOP on three baseline SFT models: Mistral-7B-Ours, MetaMath-Mistral-7B, and MathCoder-Mistral-7B. Mistral-7B-Ours is finetuned on the 34K GSM8K samples and 47K MATH samples we collected from GPT-4. MetaMath-Mistral-7B is downloaded from the MetaMath HuggingFace repository⁴. MathCoder-Mistral-7B is finetuned using the MathCodeInstruct dataset Wang et al. (2023a), downloaded from HuggingFace⁵. We collect D_{naive} and D_{SC} as described in 2.1 using problems from GSM8K and MATH. We compare 4 methods of aligning the starting SFT models: 1. Continue finetuning the starting SFT model using supervised

³https://github.com/InternLM/InternLM-Math

⁴https://huggingface.co/meta-math/MetaMath-Mistral-7B

⁵https://huggingface.co/datasets/MathLLMs/MathCodeInstruct

432 433	Method	GSM8K	MATH	OCW	hungariar	n Mathematics	SVAMP	Simuleq
434	Step-DPO	80.4	29.3	12.5	42	33.4	80.4	72.6
435	SCDPO (ours)	81.7	29.3	15.4	42	35.0	81.6	73.2

438

Table 4: Comparison between our method and Step-DPO on MetaMath-Mistral-7B.

finetuning with preferred solutions from D_{naive} (SFT-continued). 2. Doing naive DPO training with D_{naive} (SFT+DPO). 3. Doing naive DPO training with the same amount of training pairs as the SCDPO training, expanded from D_{naive} (SFT+DPO_(d-e)). 4. Doing SCDPO training with D_{naive} and D_{SC} (SFT+DPO+SC).

The results are shown in Tab. 2. The purpose of SFT+DPO_(d-e) is to rule out the possibility that the performance gain of SCDPO is the effect of more training samples. SFT-continued shows no obvious gains, likely due to the fact that the models has already been finetuned on many solutions from GSM8K and MATH. As demonstrated in Tab 2, on all three SFT baseline models, SCDPO shows superior performance compared to DPO. This can be attributed to SCDPO's more detailed supervision on the reasoning steps of the math solutions, demonstrating the effectiveness of our method.

We also compare our method with Step-DPO Lai et al. (2024), a work concurrent with ours, which
uses GPT-4 to locate erroneous steps. As Step-DPO is in the Chain-of-Thought format, we train
MetaMath-Mistral-7B using the dataset and code of Step-DPO. As shown in Tab. 4, our method outperforms Step-DPO on most datasets without relying on any stronger LLMs (e.g., GPT-4), demonstrating the effectiveness of our approach.

455 456 457

4.3 ANALYSIS OF THE INCREASING TEMPERATURE DESIGN

We present the result of using different temperature during sampling of erroneous steps in Tab. 3. 458 Originally, we tried sampling for the incorrect solutions at the same temperature as the correct 459 solutions (1.0). However, we observed that generated error steps are less diverse than we hoped. 460 Also, as shown in Fig. 4, the accuracy decreases with the increase of temperature, as the generation 461 becomes less stable. We then tried raising the temperature to 1.3, and found that a notable part of 462 the generated solutions contains incomprehensible strings at later steps due to accumulated errors. 463 Finally, we settled on raising the temperature gradually, which enables more diversity while lowering 464 the frequency of generating unintelligible sentences. As shown in Tab. 3, this method also performs 465 best in the pilot experiments.

466 467 468

5 RELATED WORK

469 LLM for Mathematical Reasoning. Prior works have explored various methods to enhance math-470 ematical reasoning abilities of LLMs. Prompting methods, such as Chain-of-Thought Wei et al. 471 (2023a), Tree-of-Thought Yao et al. (2023), PAL Gao et al. (2023), Program-of-Thought Chen 472 et al. (2023), and CSV Zhou et al. (2023), use carefully engineered prompts to bring out LLMs' 473 mathematical skills without changing their parameters. Other works optimize parameters of LLMs 474 for enhanced mathematical reasoning through either pretraining or finetuning. Llemma Azerbayev 475 et al. (2024), and MathPile Wang et al. (2023b) continue pretraining LLMs on large amounts of 476 math-related data, while RFT Yuan et al. (2023), Mammoth Yue et al. (2023), MathCoder Wang et al. (2023a), WizardMath Luo et al. (2023), ToRA Gou et al. (2024), MetaMath Yu et al. (2024), 477 MathGLM Yang et al. (2023b), and MathGenie Lu et al. (2024) finetune pretrained models on 478 question-solution pairs. These methods effectively improves LLMs' ability to solve challenging 479 mathematical problems, demonstrating impressive performance on mathematical benchmarks such 480 as GSM8K Cobbe et al. (2021), MATH Hendrycks et al. (2021), etc. Our work builds upon models 481 that have undergone pretraining and finetuning, using DPO to further enhance their mathematical 482 abilities. 483

Improving Mathematical Reasoning Using Relative Feedback. Reinforcement learning from human (or AI) feedback Christiano et al. (2023); Bai et al. (2022) as well as several direct alignment methods Rafailov et al. (2024b); Azar et al. (2023); Zhao et al. (2023); Pal et al. (2024); Ethayarajh

486 et al. (2024); Liu et al. (2024) have proven effective on various downstream tasks. Our method 487 make use of DPO Rafailov et al. (2024b), introducing a novel way to construct the DPO training 488 data for better enhancement of mathematical abilities of LLMs. Previous works using reinforcement 489 learning or direct alignment methods for improving mathematical reasoning utilize either outcome 490 supervision or process supervision. Outcome supervision such as Shao et al. (2024) is simple and use the outcome of a solution as supervision signal. Lightman et al. (2023) found that process su-491 pervision offers better performance than outcome supervision, but needs expert and detailed human 492 or AI annotation, which is difficult to acquire. Math-Shepherd Wang et al. (2023a) and Process 493 Reward Synthesizing Jiao et al. (2024) estimate process rewards with multiple decoding rationales 494 at each step, and train a reward model with the synthesized rewards. Other works such as Xie et al. 495 (2024), Yuan et al. (2024) and Chen et al. (2024) use tree structure to provide fine-grained super-496 vision, often relying on a critique model to decide the correctness of reasoning steps. Concurrent 497 works such as Setlur et al. (2024) and Lai et al. (2024) rely on GPT-4 to synthesize data. Step-498 DPO (Lai et al., 2024) uses GPT-4 for erroneous step localization, which is less cost-effective. In 499 comparison, our method uses increasing temperature to start generating erroneous steps from inter-500 mediate steps of a correct solution, and directly contrast erroneous steps with correct steps, offering 501 a simpler, more cost-effective alternative with high performance, without relying on any stronger LLMs (e.g. GPT-4). 502

503 504

505

6 LIMITATIONS AND FUTURE WORK

506 Our work contains the following limitations, and we leave them for future work. Firstly, our work is 507 conducted on purely linguistic models, which struggle to solve mathematical problems requiring an 508 understanding of images. Secondly, due to the stepwise attribute of SCDPO, it is not very effective 509 on solution formats consisting of pure code. It only works on solutions consisting of natural language 510 chain of thought or interleaved natural language and code. A method to properly enhance pure code 511 solutions needs to be derived, which we leave for future work. Thirdly, as with all language models, 512 our models can potentially generate hallucinations or produce misleading solutions, which can have a negative effect. Finally, while the data construction method distributionally narrow down the steps 513 likely to be erroneous, it does not indicate the exact step the error occurs, a problem inherent with 514 synthetic process supervision methods. Additionally, our work focuses on mathematical problem-515 solving, without discussing other reasoning tasks such as code generation, theorem proving, etc. We 516 plan to explore them in future works. 517

518 519

520

527 528

529

7 CONCLUSION

In this work, we propose Step-Controlled DPO (SCDPO), a method to automatically introduce step wise error supervision to the process of DPO training by generating dispreferred samples that start
 making errors at a specified step. SCDPO effectively enhances the mathematical reasoning abilities
 of LLMs. The 20B model trained with SCDPO on both English and Chinese data achieves high
 scores on 10 representative mathematical datasets, consistently outperforming naive DPO, demonstrating the effectiveness of our method.

- References
- Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal
 Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human
 preferences, 2023. URL https://arxiv.org/abs/2310.12036.

Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q. Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for mathematics, 2024. URL https://arxiv.org/abs/2310.10631.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli TranJohnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse,

Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna
Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario
Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai:
Harmlessness from ai feedback, 2022. URL https://arxiv.org/abs/2212.08073.

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui 550 Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye 551 Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting 552 Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li, Linyang Li, 553 Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun 554 Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang 555 Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, 556 Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong, 558 Chao Xu, Ruiliang Xu, Hang Yan, Yirong Yan, Xiaogui Yang, Haochen Ye, Huaiyuan Ying, Jia 559 Yu, Jing Yu, Yuhang Zang, Chuyu Zhang, Li Zhang, Pan Zhang, Peng Zhang, Ruijie Zhang, Shuo Zhang, Songyang Zhang, Wenjian Zhang, Wenwei Zhang, Xingcheng Zhang, Xinyue Zhang, 560 Hui Zhao, Qian Zhao, Xiaomeng Zhao, Fengzhe Zhou, Zaida Zhou, Jingming Zhuo, Yicheng 561 Zou, Xipeng Qiu, Yu Qiao, and Dahua Lin. Internlm2 technical report, 2024. URL https: 562 //arxiv.org/abs/2403.17297. 563
- Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. Step-level value preference optimization
 for mathematical reasoning, 2024. URL https://arxiv.org/abs/2406.10858.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks, 2023. URL https: //arxiv.org/abs/2211.12588.
- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
 reinforcement learning from human preferences, 2023. URL https://arxiv.org/abs/
 1706.03741.
- 573 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, 574 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John 575 Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv. 576 org/abs/2110.14168.

577

578

579

580

- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization, 2024. URL https://arxiv.org/abs/ 2402.01306.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
 Graham Neubig. Pal: Program-aided language models, 2023. URL https://arxiv.org/
 abs/2211.10435.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and
 Weizhu Chen. Tora: A tool-integrated reasoning agent for mathematical problem solving, 2024.
 URL https://arxiv.org/abs/2309.17452.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL https://arxiv.org/abs/2103.03874.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gi anna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le

613

620

624

625

626

627

630

631

636

594 Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024. URL https://arxiv.org/abs/2401.04088. 596

- Fangkai Jiao, Chengwei Qin, Zhengyuan Liu, Nancy F. Chen, and Shafiq Joty. Learning planning-597 based reasoning by trajectories collection and process reward synthesizing, 2024. URL https: 598 //arxiv.org/abs/2402.00658.
- 600 Nate Kushman, Yoav Artzi, Luke Zettlemoyer, and Regina Barzilay. Learning to automatically 601 solve algebra word problems. In Kristina Toutanova and Hua Wu (eds.), Proceedings of the 602 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 603 pp. 271–281, Baltimore, Maryland, June 2014. Association for Computational Linguistics. doi: 604 10.3115/v1/P14-1026. URL https://aclanthology.org/P14-1026.
- Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xiangru Peng, and Jiaya Jia. Step-dpo: Step-606 wise preference optimization for long-chain reasoning of llms, 2024. URL https://arxiv. 607 org/abs/2406.18629. 608
- 609 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-610 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam 611 Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with lan-612 guage models, 2022. URL https://arxiv.org/abs/2206.14858.
- Minpeng Liao, Wei Luo, Chengxi Li, Jing Wu, and Kai Fan. Mario: Math reasoning with code 614 interpreter output – a reproducible pipeline, 2024. URL https://arxiv.org/abs/2401. 615 08190. 616
- 617 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan 618 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step, 2023. URL 619 https://arxiv.org/abs/2305.20050.
- Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J. Liu, and Jialu 621 Liu. Statistical rejection sampling improves preference optimization, 2024. URL https:// 622 arxiv.org/abs/2309.06657. 623
 - Zimu Lu, Aojun Zhou, Houxing Ren, Ke Wang, Weikang Shi, Junting Pan, Mingjie Zhan, and Hongsheng Li. Mathgenie: Generating synthetic data with question back-translation for enhancing mathematical reasoning of llms, 2024. URL https://arxiv.org/abs/2402.16352.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, 628 Qingwei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathemati-629 cal reasoning for large language models via reinforced evol-instruct, 2023. URL https: //arxiv.org/abs/2308.09583.
- 632 Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Zhiqiang Hu, Chenhui Shen, 633 Yew Ken Chia, Xingxuan Li, Jianyu Wang, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue 634 Deng, Sen Yang, Chaoqun Liu, Hang Zhang, and Lidong Bing. Seallms – large language models 635 for southeast asia, 2024. URL https://arxiv.org/abs/2312.00738.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-637 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red 638 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-639 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher 640 Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-641 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, 642 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, 643 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey 644 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila 645 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, 646 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-647 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan

648 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal-649 lacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan 650 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, 651 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 652 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook 653 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel 654 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen 655 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel 656 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, 657 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv 658 Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, 659 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, 660 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 661 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-662 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel 663 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, 665 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, 666 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra 667 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, 668 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-669 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, 670 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 671 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, 672 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-673 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-674 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, 675 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-676 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming 677 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao 678 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL 679 https://arxiv.org/abs/2303.08774. 680

 Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White.
 Smaug: Fixing failure modes of preference optimisation with dpo-positive, 2024. URL https: //arxiv.org/abs/2402.13228.

684

688

689

690 691

- Keiran Paster. Testing language models on a held-out high school national finals
 exam. https://huggingface.co/datasets/keirp/hungarian_national_hs_
 finals_exam, 2023.
 - Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems?, 2021. URL https://arxiv.org/abs/2103.07191.
- Martin L. Puterman. Markov decision processes: Discrete stochastic dynamic programming. In
 Wiley Series in Probability and Statistics, 1994. URL https://api.semanticscholar.
 org/CorpusID:260799408.
- Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From r to q*: Your language model is secretly a q-function, 2024a. URL https://arxiv.org/abs/2404.12358.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024b. URL https://arxiv.org/abs/2305.18290.
- 701 David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical reasoning abilities of neural models, 2019. URL https://arxiv.org/abs/1904.01557.

- Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. Rl on incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold, 2024. URL https://arxiv.org/abs/2406.14532.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402. 03300.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Language models are multilingual chain-of-thought reasoners. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023. URL https://openreview.net/pdf?id=fR3wGCk-IXp.

715 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-716 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, 717 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy 718 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, 719 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel 720 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, 721 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, 722 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh 723 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen 724 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, 725 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 726 2023. URL https://arxiv.org/abs/2307.09288. 727

- Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi
 Song, Mingjie Zhan, and Hongsheng Li. Mathcoder: Seamless code integration in llms for en hanced mathematical reasoning, 2023a. URL https://arxiv.org/abs/2310.03731.
- Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang
 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024.
 URL https://arxiv.org/abs/2312.08935.
- Zengzhi Wang, Rui Xia, and Pengfei Liu. Generative ai for math: Part i mathpile: A billion-token scale pretraining corpus for math, 2023b. URL https://arxiv.org/abs/2312.17120.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023a. URL https://arxiv.org/abs/2201.11903.
- Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and Bin Wang. Cmath: Can your language model pass chinese elementary school math test?, 2023b. URL https://arxiv.org/abs/2306.
 16636.
- Yuxi Xie, Anirudh Goyal, Wenyue Zheng, Min-Yen Kan, Timothy P. Lillicrap, Kenji Kawaguchi, and Michael Shieh. Monte carlo tree search boosts reasoning via iterative preference learning, 2024. URL https://arxiv.org/abs/2405.00451.
- Yifan Xu, Xiao Liu, Xinghan Liu, Zhenyu Hou, Yueyan Li, Xiaohan Zhang, Zihan Wang, Aohan Zeng, Zhengxiao Du, Wenyi Zhao, Jie Tang, and Yuxiao Dong. Chatglm-math: Improving math problem-solving in large language models with a self-critique pipeline, 2024. URL https://arxiv.org/abs/2404.02893.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao,

765

769

782

784

785

790

794

796

797 798

799

800

801

756 Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng 758 Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, 759 Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024. URL 760 https://arxiv.org/abs/2407.10671.

- 762 Liu Yang, Haihua Yang, Wenjun Cheng, Lei Lin, Chenxia Li, Yifu Chen, Lunan Liu, Jianfei Pan, Tianwen Wei, Biye Li, Liang Zhao, Lijie Wang, Bo Zhu, Guoliang Li, Xuejie Wu, Xilin Luo, and 763 764 Rui Hu. Skymath: Technical report, 2023a.
- Zhen Yang, Ming Ding, Qingsong Lv, Zhihuan Jiang, Zehai He, Yuyi Guo, Jinfeng Bai, and Jie 766 Tang. Gpt can solve mathematical problems without a calculator, 2023b. URL https:// 767 arxiv.org/abs/2309.03241. 768
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik 770 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023. 771 URL https://arxiv.org/abs/2305.10601. 772
- 773 Yi. A series of large language models trained from scratch by developers at 01-ai. https:// 774 github.com/01-ai/Yi, 2023. 775
- 776 Huaiyuan Ying, Shuo Zhang, Linyang Li, Zhejian Zhou, Yunfan Shao, Zhaoye Fei, Yichuan Ma, 777 Jiawei Hong, Kuikun Liu, Ziyi Wang, Yudong Wang, Zijian Wu, Shuaibin Li, Fengzhe Zhou, 778 Hongwei Liu, Songyang Zhang, Wenwei Zhang, Hang Yan, Xipeng Qiu, Jiayu Wang, Kai Chen, 779 and Dahua Lin. InternIm-math: Open math large language models toward verifiable reasoning, 780 2024. URL https://arxiv.org/abs/2402.06332. 781
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhen-783 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models, 2024. URL https://arxiv.org/abs/2309.12284.
- 786 Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin 787 Chen, Ruobing Xie, Yankai Lin, Zhenghao Liu, Bowen Zhou, Hao Peng, Zhiyuan Liu, and Maosong Sun. Advancing llm reasoning generalists with preference trees, 2024. URL https: 788 //arxiv.org/abs/2404.02078. 789
- Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuangi Tan, Chang Zhou, 791 and Jingren Zhou. Scaling relationship on learning mathematical reasoning with large language 792 models, 2023. URL https://arxiv.org/abs/2308.01825. 793
 - Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mammoth: Building math generalist models through hybrid instruction tuning, 2023. URL https://arxiv.org/abs/2309.05653.
 - Wei Zhao, Mingyue Shang, Yang Liu, Liang Wang, and Jingming Liu. Ape210k: A large-scale and template-rich dataset of math word problems, 2020. URL https://arxiv.org/abs/ 2009.11506.
- 802 Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu. Slic-hf: 803 Sequence likelihood calibration with human feedback, 2023. URL https://arxiv.org/ 804 abs/2305.10425. 805
- 806 Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, 807 Linqi Song, Mingjie Zhan, and Hongsheng Li. Solving challenging math word problems using 808 gpt-4 code interpreter with code-based self-verification, 2023. URL https://arxiv.org/ 809 abs/2308.07921.

810 811	Data	GSM8K	MATH	
812	D _{SC}	79.0%	46.2%	
813	$D_{\text{naive}} + D_{\text{SC}}$	80.1%	47.7%	
814				

Table 5: Ablation study of using and not using D_{naive} during training. The starting SFT model is Mistral-7B-Ours.

816 817 818

819 820

821

822

823

824

825 826

827 828

829

830

815

A CREDIT ASSIGNMENT ANALYSIS EXAMPLES

In this section, we present several other credit assignment analysis examples, comparing SCDPO to DPO. Fig. 6, Fig. 7 and Fig. 8 show examples of part of the solutions of questions taken from GSM8K and MATH datasets, colored with the DPO implicit reward of each token (darker is higher). As demonstrated in the examples, SCDPO is better than DPO at identifying the errors in the reasoning steps.

B Ablation Study of D_{NAIVE} and D_{SC}

To demonstrate the necessity of combining D_{naive} and D_{SC} , we conduct experiment of using only D_{SC} in DPO training. The results are presented in Tab. 5. As demonstrated in the table, combining D_{naive} and D_{SC} results in better performance than only using D_{SC} during DPO training. This is likely because D_{naive} helps regulate the general format of the generated solutions.

831 832 833

834

840 841

842

843

844

845

846

847

848 849

850

C ANALYSIS OF GENERATED ERRORS

In this section, we provide quantitative analysis of the erroneous steps generated. We observed seven main kinds of errors: value misusage, condition misinterpretation, coding error, commonsense error, math concept or understanding error, math calculation error, unintelligible strings. The errors are explained as follows:

- Value misusage: misusing values in places where another value should have been used.
- Condition misinterpretation: incorrectly interpreting the meaning or indications of conditions.
- Coding error: making mistakes in code snippets that causes errors.
- Common sense error: misunderstanding of common sense.
- Math concept or understanding error: incorrect recollection or understanding of math concepts.
 - Math calculation error: mistakes when making mathematical calculations.
 - Unintelligible strings: generation of unintelligible strings that does not represent meaningful reasoning errors.

We randomly sampled 100 incorrect solutions in the training data of SCDPO, and counted the number of each type of error. The result is presented in Fig. 5. As demonstrated in the chart, the reasoning errors generated is diverse, distributed evenly among the different types. Only 4% of the incorrect solutions contain unintelligible strings, demonstrating that the design of gradually increasing temperature can mostly avoid the occurance of meaningless errors.

D ERROR RATE OF INTERMEDIATE STEPS WHEN THE FINAL ANSWER IS CORRECT

859 860

857

858

861 In this section, we discuss the error rate of intermediate steps in solutions that reaches the correct 862 final answer. As we mentioned in the main paper, we randomly sampled 87 solutions that reach 863 correct final answer, and of the 369 reasoning steps in these solutions, only 2 contain errors, which 863 is a very small percentage (about 0.5%). The 2 erroneous steps are in a question whose answer is to



at the beginning of the generation of erroneous steps, most of them occur early in the generation, as

918	# Let's denote the number of pencils Antonio has as x.	# Let's denote the number of pencils Antonio has as x.
919	# According to the problem, Mitchell has x + 6 pencils.	# According to the problem, Mitchell has x + 6 pencils.
920	# Given that Mitchell has 30 pencils, we can set up the equation:	# Given that Mitchell has 30 pencils, we can set up the equation:
921		
922	# v t (v t 6) - 30	$+ \times + (\times + 6) = 30$
923		
924	H Calification for a single set	H Onlying for a single set
925	# Solving for x gives us:	# Solving for X gives us:
000	from sympy import symbols, Eq, solve	from sympy import symbols, Eq, solve

Figure 7: Credit assignment of part of a solution for a GSM8K problem. Each token is colored corresponding to the DPO implicit reward as expressed in Eq. 6 (darker is higher). The left is the credit assignment of SCDPO, which correctly highlighted the error – Mitchell has 30 pencils, and Antonio has 6 less pencils than Michell, which is 30 - 6, so the introduction of x is not needed, and x + (x + 6) = 30 is incorrect, while the credit assignment of DPO on the right fails to highlight it.

To solve this problem, we need to find \(n \) such that the sum of	To solve this problem, we need to find \(n \) such that the sum of
an arithmetic series satisfies the given congruence.	an arithmetic series satisfies the given congruence.
The arithmetic series starts at 1 and has a common difference of	The arithmetic series starts at 1 and has a common difference of
5. The last term, 101, can be written as \(5k \), where \(k = 20 \)	5. The last term, 101, can be written as $(5k)$, where $(k = 20)$.

Figure 8: Credit assignment of part of a solution for a MATH problem. Each token is colored corresponding to the DPO implicit reward as expressed in Eq. 6 (darker is higher). The left is the credit assignment of SCDPO, which correctly highlighted the error -101 cannot be written as 5k where k = 20, while the credit assignment of DPO on the right fails to highlight the error.

the earlier steps often make decisions that have a large influence on the correctness of the solution.
We present the distribution of the step index where the first actual error occurs in the 100 SCDPO
rejected solutions, as well as the distribution of the index of the step where generation of erroneous
steps begins. As shown in Tab. 6, the two distributions are closely related, demonstrating that the
index of the starting point of error generation have a strong effect on the step index of the actual
error. The step index of the actual error is often close to the starting point of the error generation.

We also manipulate the distribution of the step index where error generation begins. Specifically, we train the Mistral-7B-SFT model with SCDPO data where the starting index was either less than or equal to 4, or greater than 4. As shown in Tab. 7, limiting the error starting point index decreases performance compared to not imposing any such limitation.

Step Index	2	3	4	5	6	7	8	9	10	11	12
Actual Error Distr. Starting Point Distr.	17 32.8	21 23.1	20 15.3	12 11.4	8 6.97	9 4.72	5 2.53	3 1.70	0 0.77	4 0.44	1 0.16
Table 6: The distributiorejected solutions and trejected solutions.	n of the	e actu ibutic	al erro on of t	or star he sta	ting s arting	step in point	ndex in of the	n 100 e erro	rando r gene	omly seration	sampl n of ti
Data	GSM8	<u> </u>	IATH	00	CW	Huns	garian	Matl	nemati	cs SV	VAMP
Data SFT+DPO+SC	GSM8 80.1	KN	IATH 47.7	00	CW 2.4	Hung	garian	Matl	nemati	cs SV	VAMP 82.3
Data SFT+DPO+SC SFT+DPO+SC ($k \le 4$)	GSM8 80.1 80.7	KM	IATH 47.7 46.7	OC 22 17	C W 2.4 7.7	Hung 61 53	garian 1 3	Matl	nemati 6.5 4.7	cs SV	VAMP 82.3 82.0

Table 7: Comparison of performance on various datasets using different ranges of error generation starting point during training. "k" is the index of the starting step of error generation. Limiting the starting index slightly affects performance compared to using the full range of indices.