

# WizardMath: EMPOWERING MATHEMATICAL REASONING FOR LARGE LANGUAGE MODELS VIA *Reinforced Evol-Instruct*

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models (LLMs), such as GPT-4, have shown remarkable performance in natural language processing (NLP) tasks, including challenging mathematical reasoning. However, most existing open-source models are only pre-trained on large-scale internet data and without math-related optimization. In this paper, we present *WizardMath*, which enhances the mathematical CoT reasoning abilities of LLMs without using external python tools, by applying our proposed *Reinforcement Learning from Evol-Instruct Feedback (RLEIF)* method to the domain of math. Through extensive experiments on two mathematical reasoning benchmarks, namely GSM8k and MATH, we reveal the extraordinary capabilities of our model. Remarkably, *WizardMath-Mistral 7B* surpasses top-tier open-source LLMs by a substantial margin with higher data efficiency. Furthermore, *WizardMath 70B* even outperforms GPT-3.5-Turbo, Claude 2, Gemini Pro and GPT-4-early-version. Additionally, our preliminary exploration highlights the pivotal role of instruction evolution and process supervision in achieving exceptional math performance.

## 1 INTRODUCTION

Recently, Large-scale language models (LLMs) have garnered significant attention and become the go-to approach for numerous natural language processing (NLP) tasks, including open domain conversation (Ouyang et al., 2022; OpenAI, 2023; Touvron et al., 2023a), coding (Chen et al., 2021; Wang et al., 2021; Li et al., 2023b) and math (Taylor et al., 2022; Lewkowycz et al., 2022; Shao et al., 2024; Yang et al., 2024). A conspicuous example is ChatGPT<sup>1</sup>, developed by OpenAI. This model uses extensive pre-training on large-scale internet data and further fine-tuning with specific instruction data and methods. As a result, it achieves state-of-the-art zero-shot performance on various benchmarks. Subsequently, Anthropic, Google, and Meta also launched their competitive products one after another. Notably, Meta’s series of Llama (Touvron et al., 2023a;b; Dubey et al., 2024) have sparked an open-source revolution and quickly narrowed the gap with those closed-source LLMs. This trend also gradually stimulates the releases of Mistral (Jiang et al., 2023), Alpaca (Taori et al., 2023), Vicuna (Chiang et al., 2023), and WizardLM (Xu et al., 2023), etc. However, these open models still struggle with the scenarios which require complex multi-step quantitative reasoning, such as solving mathematical and science challenges (Ahn et al., 2024; Long et al., 2024).

Chain-of-thought (CoT) (Wei et al., 2022) proposes to design better prompts to generate step-by-step solutions, which can lead to improved performance. Self-Consistency (Wang et al., 2022) also achieves remarkable performance on many reasoning benchmarks, which generates several possible answers from the model and selects the correct one based on majority vote (Fu et al., 2022). Llemma (Azerbayev et al., 2023) and MathPile (Wang et al., 2023f) continue pretraining LLMs with math corpus to improve domain capacity. MetaMath (Yu et al., 2023b) and Xwin-Math (Li et al., 2024a) bootstraps mathematical questions by augmenting the question from multiple perspectives. MAMmoTH (Yue et al., 2023) and TORA (Gou et al., 2023) presents a unique hybrid of CoT and program-of-thought (PoT) to ensure extensive coverage of diverse fields in math. Recently, Evol-Instruct is an effective method for large-scale data synthesis using LLMs. It has been widely verified and proven to be effective in enhancing the model’s instruction following capability. It employs

<sup>1</sup> <https://openai.com/>

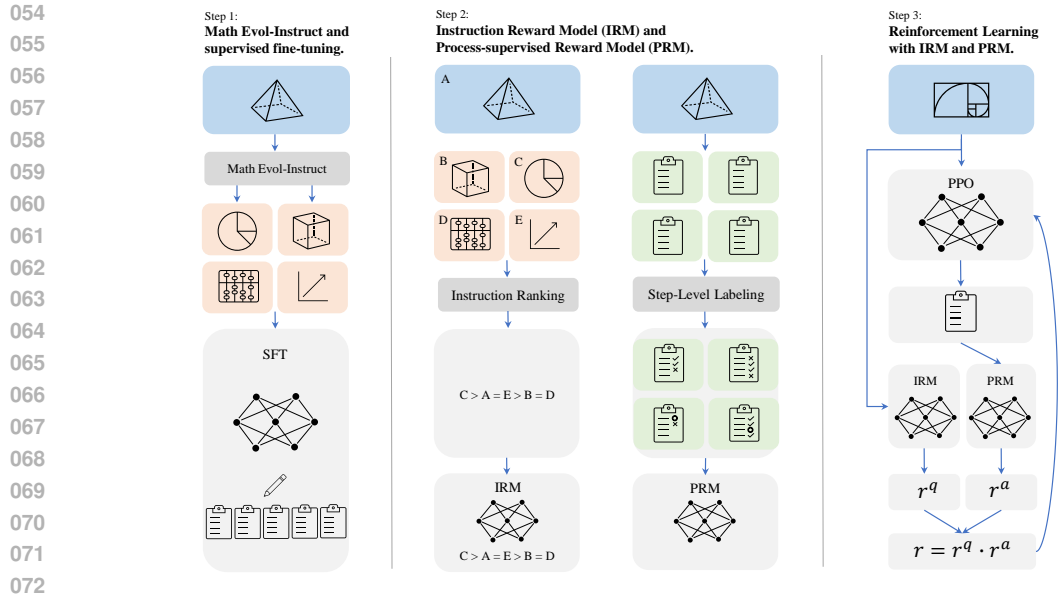


Figure 1: A diagram illustrating the three steps of our *Reinforcement Learning from Evol-Instruct Feedback (RLEIF)*

In-depth Evolving and In-breadth Evolving to automate the generation of diverse and complex open-domain instructions using LLMs, instead of relying on human-crafted instruction datasets. In-depth Evolving incrementally enhances instruction complexity by introducing additional constraints, deepening, concretizing, increasing reasoning steps, and complicating input. In-breadth Evolving focuses on improving topic diversity and dataset richness by creating entirely new instructions. To enhance the correctness of each step in the model’s generation process, (Wang et al., 2024a; Chen et al., 2024a; Lightman et al., 2023) finds that process supervision with reinforcement learning significantly outperforms outcome supervision for solving challenging MATH problems.

Inspired by *Evol-Instruct* and Process-supervised Reinforcement Learning, this work aims to enhance the mathematical reasoning abilities of the LLMs. As shown in the Figure 1, we propose a new method named *Reinforcement Learning from Evol-Instruct Feedback (RLEIF)*, which could firstly generate diverse math instructions data by brand-new *Math Evol-Instruct*, which includes two downward evolution and upward evolution progress to produce the grade school math and challenging high school math respectively. However different from WizardLM (Xu et al., 2023) and WizardCoder (Luo et al., 2023), which mainly focus on the SFT stage and are susceptible to learning hallucinated information from the teacher model, we innovatively introduce PRM to address the False-Positive issue in the problem-solving process. Moreover, to prevent instruction evolution from spiraling out of control, we incorporate an instruction reward model (IRM) as a mitigating strategy. Thus, we train an instruction reward model (IRM) and a process-supervised reward model (PRM) (Lightman et al., 2023; Uesato et al., 2022; Wang et al., 2024a; Chen et al., 2024a), the former indicates the quality of the evolved instruction and the latter offers feedback for each reasoning step in the solution. Initially, we finetune LLMs with the evolved math data. Immediately, we leverage GPT-4 to produce the ranking order of instructions, and the correctness of each reasoning step, then optimize the LLMs to obtain the reward models. Finally, we implement the step-by-step PPO to train our *WizardMath*.

We perform experiments on two widely used mathematical reasoning benchmarks, namely GSM8k (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) covering math problems from grade to high school levels, the results show that our *WizardMath* outperforms all other open-source LLMs at the same model size, achieving state-of-the-art performance. For instance, *WizardMath-70B* significantly outperforms MetaMath-70B by a significant margin on GSM8k (92.8 vs. 82.3) and on MATH (58.6 vs. 26.6). Specifically, *WizardMath-Mistral-7B* observed a substantial improvement in pass@1 with an increase of +12.8 (90.7 vs. 77.9) on GSM8k, and +26.8 (55.4 vs. 28.6) on MATH compared to MetaMath-Mistral-7B. Notably, our 70B model even also significantly surpasses those powerful

proprietary LLMs, such as GPT-3.5-Turbo, Claude 2 (Bai et al., 2022), Mistral Medium (Jiang et al., 2024a), Gemini-Pro (Team, 2023), PaLM-2 (Anil et al., 2023) and GPT-4-early-version.

The main contributions of this work are as follows:

- We introduce *WizardMath* model, which enhances the LLMs’ mathematical reasoning abilities across a range of problem difficulties, from grade to high school levels.
- We propose a new fully AI-powered automatic reinforcement learning method, *Reinforcement Learning from Evol-Instruct Feedback (RLEIF)*, alongside *Math Evol-Instruct* and Process Supervision, for improving reasoning performance.
- *WizardMath* surpasses top-tier open-source LLMs by a substantial margin with higher data efficiency and also significantly outperforms various proprietary LLMs on both GSM8k and MATH, demonstrate the effectiveness of our *RLEIF*.

## 2 RELATED WORK

**Large Language Models.** LLMs have significantly advanced Natural Language Processing, with models like OpenAI’s GPT Series (Brown et al., 2020a; OpenAI, 2023), Anthropic’s Claude (Bai et al., 2022), Google’s PaLM (Chowdhery et al., 2022; Anil et al., 2023), Gemini (Team, 2023), and Gemma (Team et al., 2024) featuring billions of parameters and trained on massive textual datasets. The AI field has also seen a rise in open-source LLMs such as Mistral (Jiang et al., 2023), Llama Series (Touvron et al., 2023a;b; Dubey et al., 2024; Taylor et al., 2022), DeepSeek (Bi et al., 2024; Shao et al., 2024), Qwen (Bai et al., 2023; Yang et al., 2024) etc. Notably, Llama serves as a foundational model for supervised fine-tuning, leading to the development of models like Alpaca, Vicuna (Taori et al., 2023; Chiang et al., 2023).

**Large Language Models For Mathematical reasoning.** NLP models face challenges with complex reasoning, including mathematical (Long et al., 2024; Zhang et al., 2024c; Xia et al., 2024), common-sense (Talmor et al., 2019). Significant research focuses on Mathematical Word Problems (MWP), which demand understanding of mathematical concepts and multi-step reasoning (Zheng et al., 2023; Zhao et al., 2023; Yuan et al., 2023a). Models are tested on various MWP benchmarks (Roy & Roth, 2015; Hendrycks et al., 2021). Techniques like Chain-of-Thought Prompting (Wei et al., 2022), Least-to-Most prompting (Zhou et al., 2022), and Complex CoT (Fu et al., 2022) enhance reasoning by introducing multiple steps and breaking problems into sub-problems. There are some models aimed at improving math CoT reasoning skills such as MetaMath (Yu et al., 2023b), MathScale (Tang et al., 2024), Xwin-Math (Li et al., 2024a), DART-Math (Tong et al., 2024) etc. Some models enhance mathematical reasoning by integrating python tools, such as TORA (Gou et al., 2023), MAMmoTH (Yue et al., 2023), Openmathinstruct (Toshniwal et al., 2024), NuminaMath (Li et al., 2024c) etc. In our work, we mainly improve the CoT reasoning ability of mathematics without using external Python tools.

**Reinforcement Learning for Large Language Models.** State-of-the-art models often display logical errors and illusions, particularly in domains requiring complex, multi-step reasoning, leading to significant challenges (Bubeck et al., 2023; Maynez et al., 2020). Strategies such as training reward models help discriminate between desirable and undesirable outputs (Lightman et al., 2023; Wu et al., 2023b; Chen et al., 2024b). Historically, outcome-based approaches focused on algorithmic tasks (Li et al., 2016; Cai et al., 2017; Yu et al., 2023a), while recent research demonstrates the efficacy of reward models or validators in enhancing model performance (Cobbe et al., 2021; Wang et al., 2023c;d; Li et al., 2022a). Reward models have also been incorporated into reinforcement learning pipelines and employed in rejection sampling to align Large Language Models (LLMs) with human preferences (Shen et al., 2021; Bai et al., 2022; Yuan et al., 2023c; Dong et al., 2023; Song et al., 2023; Touvron et al., 2023b; Rafailov et al., 2024; Meng et al., 2024). A contrast is drawn between outcome-supervised and process-supervised reward models, with the latter being more effective at addressing discrepancies arising from incorrect reasoning paths leading to correct outcomes (Uesato et al., 2022; Zelikman et al., 2022; Creswell et al., 2022). Recent advances have promoted process-based supervision through manual annotation, significantly benefiting LLMs over outcome-based approaches (Lightman et al., 2023; Wang et al., 2024a; Sun et al., 2024; Chen et al., 2024a; Wang et al., 2024b; Zhang et al., 2024a). In this paper, we leverage AI models like ChatGPT to automatically offer process annotation to improve the efficiency of this research line.

### 3 METHOD

In this section, we elaborate on the details of our *WizardMath*. Following WizardLM and PRMs (Lightman et al., 2023), we propose *Reinforcement Learning from Evol-Instruct Feedback (RLEIF)* method, which integrates the math *Evol-Instruct* and reinforced instruction and process supervision to evolve GSM8k and MATH, and fine-tune the pre-trained language models with the evolved data and reward models.

#### 3.1 MATH EVOL-INSTRUCT

Motivated by the Evol-Instruct (Xu et al., 2023) method proposed by WizardLM and its effective application on WizardCoder (Luo et al., 2023), this work attempts to make math instructions with various complexities and diversity to enhance the pre-trained LLMs. Specifically, we adapt Evol-Instruct to a new paradigm including two evolution lines:

- 1) Downward evolution: It enhances instructions by making the questions easier. For example i): revising high difficulty questions to lower difficulty, or ii) producing a new and easier question with another different topic.
- 2) Upward evolution: Derived from original Evol-Instruct method, it deepens and generates new and harder questions by i) adding more constraints, ii) concretizing, iii) increasing reasoning.

The complete prompts of above evolution are shown in Appendix A.1. For each instruction, we use GPT-4 to evolve 5 rounds (2 downward and 3 upward) of new instructions progressively, each new one is generated by the previous round of evolution.

#### 3.2 REWARD MODELS

Considering the necessity of quality control for evolved instructions and inspired by PRMs (Lightman et al., 2023), we train two reward models to predict the quality of the instructions and the correctness of each step in the answer respectively:

**Instruction Reward Model (IRM)** This model aims to judge the quality of the evolved instructions on two aspects: i) Difficulty, and ii) Definition. To produce the ranking list training data of IRM, we leverage GPT-4 to rank the quality between those evolved instructions and original instruction. The one with high difficulty and clear definition will deserve a higher ranking. The detailed prompt of above ranking process is shown in the Appendix A.2.

Specifically, given an math instructions  $q$ , IRM ( $Q \rightarrow \mathbb{R}$ ) assigns a score to  $q$  to indicate its quality. We optimize ORM via the following pairwise ranking loss:

$$\mathcal{L}_{IRM} = -\log \sigma(r_j^q - r_k^q - m) \quad (1)$$

where  $r_j^q$  is the reward of chosen instruction and  $r_k^q$  is the reward of rejected instruction,  $m$  is the margin.

**Process-supervised Reward Model (PRM)** As there is no simple way to support highly precise process supervision without professional and expensive human-labelers, we depend on GPT-4 to provide process supervision, and ask it to assess the correctness of each step in the solutions generated by our model to produce PRM training data. The detailed prompt of above step level labeling process is shown in the Appendix A.3.

For exactly, given an math instructions  $q$  and its answer  $a$ , PRM ( $Q \times A \rightarrow \mathbb{R}^+$ ) assigns a score to each step of  $a$ , we train PRM with the following cross-entropy loss:

$$\mathcal{L}_{PRM} = \sum_{i=1}^L y_i \log r_i^a + (1 - y_i) \log(1 - r_i^a) \quad (2)$$

where  $L$  is the reasoning steps of answer  $a$ .  $y_i$  is the ground-truth label of the  $i$ -th step of answer  $a$ ,  $y_i = 1$  if  $a_i$  is correct, otherwise  $y_i = 0$ .  $r_i^a$  is the reward score (assigned by PRM) of the  $i$ -th step of answer  $a$ .

### 3.3 REINFORCEMENT LEARNING WITH IRM AND PRM

Immediately, we exploit reinforcement learning to optimize LLMs. Following (Lightman et al., 2023), we employ step by step Proximal Policy Optimization (PPO) to reward both instruction and each reasoning step.

For each math instruction  $q$  and generated answer  $a$ , we use IRM to assign instruction reward  $r^q$ , and use the minimum score across all reasoning steps to represent the final reward score  $r^a$  of the answer  $a$  assigned by PRM. Then we apply a product as the final reward of this instruction-answer pair:

$$r = r^q \cdot r^a \quad (3)$$

### 3.4 PRM FOR VERIFICATION

Following (Lightman et al., 2023) and (Li et al., 2023c), we leverage both majority voting and PRM verifier to aggregate the predictions of different reasoning paths.

$$\hat{a} = \arg \max_a \sum_{i=1}^N \mathbb{I}_{a_i=a} \cdot PRM(q, a_i) \quad (4)$$

where  $PRM(q, a_i)$  is the score of the  $i$ -th reasoning path assigned by PRM for instruction  $q$ .  $\mathbb{I}_{a_i=a}$  is an indicator function that returns 1 (or 0) if  $a_i = a$ .

## 4 EXPERIMENT

This section provides a comprehensive overview of the baseline models. Subsequently, we mainly elucidate the performance metrics of our models on two prevalent mathematical benchmarks from grade to high school problems: GSM8k (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021).

### 4.1 BASELINES

Our proposed method undergoes comparison with: (i) proprietary models, including GPT-4, GPT-3.5-Turbo, PaLM 2, Gemini, Claude, Mistral Medium and Minerva; (ii) publicly available models such as MPT, GAL, Llama-2, Mistral, and Qwen; (iii) Rejection Sampling Fine-Tuning models, which generates and aggregates accurate reasoning paths as augmented data for fine-tuning, as seen in RFT and MuggleMATH; (iv) MAMmoTH which combines CoT and PoT; (v) MetaMath, employing a bootstrapping approach to rewrite mathematical questions from multiple perspectives.

### 4.2 EXPERIMENTAL SETUP

**SFT Training Data.** Firstly, use the GSM8k and MATH training sets as the initial seed collection, then employ both upward and downward math Evol-Instruct approach for five rounds. Each round need to evolve the initial instructions 6 times, and the temperature parameter is set to 0.7. Next, we remove duplicate instructions 17k. Hence, a total of 448k unique instructions were obtained. Subsequently, 30k data were excluded by the data filtering method to avoid contamination, ultimately leaving 418k data. Finally, we use GPT-4-0613 to generate the answer with a step-by-step format, and leverage them for supervised fine-tuning.

**Reward Models Training Data.** To train the reward models, We conducted additional 5 rounds of evolution on the initial instruction set and obtain 90k instructions. we use GPT-4-0613 to rank each instruction list with the quality from 1 to 6 as the training data of IRM. To obtain the training data of PRM, We use our Llama-2 70B SFT model to generate 5 answers for each instruction, and GPT-4-0613 is employed to assign correctness judgement for each reasoning step.

**Implementation Details.** We employ our method on two open-source foundational models Llama 2 (Touvron et al., 2023b) and Mistral-7B (Jiang et al., 2023). Llama 2 encompasses three distinct parameter sizes: 7B, 13B, and 70B. We utilize GPT-4-0613 for instruction evolution and the training data construction of reward models. For SFT, we train 3 epochs, and the learning rate is  $2e-5$ ,  $1e-5$  and  $5e-6$  for Llama 2 7B/13B, 70B and Mistral-7B. The batch size is 512, and the sequence length is 2048. For the reward model, we train Llama 2 and Mistral-7B with learning rate  $4e-6$  and  $1e-6$  for one epoch. For RL, the lr is  $4e-7$  and  $1e-7$  for Llama 2 and Mistral-7B and train one epoch.

## 4.3 MAIN RESULTS

Table 1 shows the CoT (Wei et al., 2022) pass@1 results of the current state-of-the-art models on GSM8k and MATH. In this study, to ensure equitable and cohesive evaluations, we report the scores of all models within the settings of **greedy decoding and CoT without using any external python tool**.

**Comparing with the proprietary Models.**

As shown in the Table 1, our *WizardMath* demonstrates notable superiority over various proprietary LLMs on the GSM8k and MATH benchmarks in terms of pass@1:

1) *WizardMath-Llama 70B*, the largest model, demonstrated exceptional performance on the GSM8k and MATH, surpassing earlier versions of GPT-4, Claude-2, and Gemini Pro, and performing on par with GPT-4-0314. It significantly outperformed GPT-3.5-Turbo by 11.2% on GSM8k and by 15.5% on MATH.

2) *WizardMath-Mistral 7B*, the smaller-sized model, outperformed Baichuan 3 on GSM8k (90.7 vs. 87.6) and surpassed GPT-4-0314 on MATH (55.4 vs. 52.6), significantly exceeding the performance of GPT-3.5-Turbo and Gemini Pro. Meanwhile, *WizardMath-Mathstral*, trained on *Mathstral-7B-v0.1*, demonstrated performance comparable to GPT-4-turbo-0125. Additionally, *WizardMath-Qwen*, trained on *Qwen2.5-Math*, surpassed GPT-4-2024-0513 on MATH (77.8 vs. 76.6).

**Comparing with the Open-Source Models.**

The results presented in Table 1 unequivocally indicate that our *WizardMath-Llama 70B* exhibits a significant performance superiority over strong models in both the GSM8k and MATH benchmarks with higher data efficiency across the range from 0.1B to 70B parameters. The detailed results are as follows:

1) With the same model parameter size, our model surpasses the previous best model such as *MetaMath*, *MAMmoTH2-Plus*, *Xwin-Math*. Particularly, *WizardMath-Llama 70B* achieves a substantial improvement of 10.5% on GSM8K and 32.0% on MATH compared to *MetaMath-Llama 70B* in testing accuracy. In the Table 2, we show the detailed results of MATH subtopics with our *WizardMath 70B* model. Specifically, *WizardMath-Mistral 7B* also surpasses top-tier open source models, outperforming *MetaMath-Mistral 7B* with a notable margin (90.7 vs 77.9 on GSM8k) and (55.4 vs 28.6 on MATH). It demonstrates the effectiveness

Table 1: The models’ CoT pass@1 results on GSM8k and MATH without using any external python tool.

Model	Base	Params	GSM8k	MATH
Proprietary models				
GPT-o1 (OpenAI, 2023)	-	-	-	94.8
GPT-o1-mini	-	-	-	90.0
Gemini-1.5 002	-	-	-	86.5
Claude 3.5 Sonnet (Bai et al., 2022)	-	-	96.4	71.1
GPT-4o-2024-0513	-	-	96.1	76.6
GPT-4-turbo-0125 (OpenAI, 2023)	-	-	94.2	64.5
GPT-4-0314	-	-	94.7	52.6
GPT-4 (original version)	-	-	92.0	42.5
Baichuan-3 (Yang et al., 2023)	-	-	88.2	49.2
GLM-4 (GLM et al., 2024)	-	-	87.6	47.9
Gemini Pro (Team, 2023)	-	-	86.5	32.6
Claude2	-	-	85.2	32.5
GPT-3.5-Turbo	-	-	81.6	43.1
PaLM2 (Anil et al., 2023)	-	-	80.7	34.3
Minerva (Lewkowycz et al., 2022)	-	540B	58.8	33.6
GPT3.5 (Brown et al., 2020a)	-	-	57.1	-
Open-Source Models (0.1B-3B)				
GPT-2-Small (Brown et al., 2020b)	-	0.1B	6.9	5.4
GPT-2-Medium (Brown et al., 2020b)	-	0.3B	11.2	6.2
GPT-2-Large (Brown et al., 2020b)	-	0.7B	13.6	6.4
GPT-2-XL (Brown et al., 2020b)	-	1.5B	15.4	6.9
WizardMath-GPT	GPT-2-Small	0.1B	26.4	12.3
WizardMath-GPT	GPT-2-Medium	0.3B	38.7	15.6
WizardMath-GPT	GPT-2-Large	0.7B	50.1	21.2
WizardMath-GPT	GPT-2-XL	1.5B	58.9	25.4
WizardMath-Qwen	Qwen-Math-2.5	1.5B	86.7	68.6
Llama-3.2-Instruct (Dubey et al., 2024)	Llama 3.2	1B	44.4	30.6
WizardMath-Llama	Llama 3.2	1B	63.3	33.5
Llama-3.2-Instruct	Llama 3.2	3B	77.7	48.0
WizardMath-Llama	Llama 3.2	3B	85.5	49.9
Open-Source Models (7B-8B)				
Llama-2 (Touvron et al., 2023b)	-	7B	14.6	2.5
MAMmoTH-CoT (Yue et al., 2023)	Llama-2	7B	50.5	10.4
MathScale (Tang et al., 2024)	Llama-2	7B	66.3	31.1
MetaMath (Yu et al., 2023b)	Llama-2	7B	66.5	19.8
MuggleMath (Li et al., 2023a)	Llama-2	7B	68.4	-
Skywork-Math (Zeng et al., 2024)	Llama-2	7B	72.9	47.7
Math-Shepherd (Wang et al., 2024a)	Llama-2	7B	73.2	21.6
Xwin-Math (Li et al., 2024a)	Llama-2	7B	82.6	40.6
WizardMath-Llama	Llama-2	7B	84.1	43.5
Mistral-v0.1 (Jiang et al., 2023)	-	7B	42.9	12.9
MathScale (Tang et al., 2024)	Mistral-v0.1	7B	74.8	35.2
MMIQc (Liu & Yao, 2024)	Mistral-v0.1	7B	74.8	36.0
MetaMath (Yu et al., 2023b)	Mistral-v0.1	7B	77.9	28.6
KPMATH-Plus (Huang et al., 2024b)	Mistral-v0.1	7B	82.1	46.8
DART-Math (Tong et al., 2024)	Mistral-v0.1	7B	82.6	43.5
Skywork-Math (Zeng et al., 2024)	Mistral-v0.1	7B	83.9	51.2
Math-Shepherd (Wang et al., 2024a)	Mistral-v0.1	7B	84.1	33.0
MAMmoTH2-Plus (Yue et al., 2024)	Mistral-v0.1	7B	84.7	45.0
JiuZhang3.0 (Zhou et al., 2024)	Mistral-v0.1	7B	88.6	52.8
Xwin-Math (Li et al., 2024a)	Mistral-v0.1	7B	89.2	43.7
WizardMath-Mistral	Mistral-v0.1	7B	90.7	55.4
WizardMath-Mistral	Mistral-v0.3	7B	90.4	55.6
WizardMath-Mathstral	Mathstral-v0.1	7B	93.8	70.9
WizardMath-Qwen	Qwen2.5-Math	7B	93.9	77.8
WizardMath-Qwen	Qwen2.5	7B	94.0	74.5
DeepSeekMath-Base (Shao et al., 2024)	-	7B	64.2	36.2
NuminaMath-CoT (Li et al., 2024c)	DeepSeekMath	7B	75.4	55.2
MMIQc (Liu & Yao, 2024)	DeepSeekMath	7B	79.0	45.3
KPMATH-Plus (Huang et al., 2024b)	DeepSeekMath	7B	83.9	48.8
DeepSeekMath-RL (Shao et al., 2024)	DeepSeekMath	7B	88.2	51.7
DART-Math (Tong et al., 2024)	DeepSeekMath	7B	88.2	52.9
WizardMath-DeepSeek	DeepSeekMath	7B	91.0	64.6
MetaMath (Yu et al., 2023b)	Llama 3	8B	77.3	20.6
MMIQc (Liu & Yao, 2024)	Llama 3	8B	77.6	29.5
DART-Math (Tong et al., 2024)	Llama 3	8B	82.5	45.3
MAMmoTH2-Plus (Yue et al., 2024)	Llama 3	8B	84.1	42.8
Llama 3.1-Instruct (Dubey et al., 2024)	Llama 3	8B	84.5	51.9
JiuZhang3.0 (Zhou et al., 2024)	Llama 3	8B	88.6	51.0
WizardMath-Llama	Llama 3	8B	90.3	58.8
Open-Source Models (13B)				
Llama-2 (Touvron et al., 2023b)	-	13B	28.7	3.9
MAMmoTH-CoT (Yue et al., 2023)	Llama 2	13B	56.3	12.9
MathScale (Tang et al., 2024)	Llama 2	13B	71.3	33.8
MetaMath (Yu et al., 2023b)	Llama 2	13B	72.3	22.4
MuggleMath (Li et al., 2023a)	Llama 2	13B	74.0	-
KPMATH-Plus (Huang et al., 2024b)	Llama 2	13B	81.6	41.0
Xwin-Math (Li et al., 2024a)	Llama 2	13B	88.1	44.9
WizardMath-Llama	Llama 2	13B	89.7	50.6
Open-Source Models (70B)				
Llama-2 (Touvron et al., 2023b)	-	70B	56.8	13.5
MAMmoTH-CoT (Yue et al., 2023)	Llama-2	70B	72.4	21.1
MetaMath (Yu et al., 2023b)	Llama-2	70B	82.3	26.6
KPMATH-Plus (Huang et al., 2024b)	Llama-2	70B	87.4	48.6
Xwin-Math (Li et al., 2024a)	Llama-2	70B	90.6	52.8
WizardMath-Llama	Llama-2	70B	92.8	58.6

Table 2: Results of pass@1 (%) on MATH subtopics (i.e., Intermediate Algebra, Geometry) with WizardMath 70B model.

MATH subtopics	WizardMath 70B
Intermediate Algebra	36.3
Precalculus	38.9
Geometry	48.3
Number Theory	58.5
Counting & Probability	54.8
Prealgebra	74.6
Algebra	78.5
Overall	<b>58.6</b>

Table 3: Explore the effects of PRM and IRM during PPO training.

Models	GSM8K	MATH
GPT-2-XL-1.5B: WizardMath-SFT	51.9	18.3
+ PRM	55.8	22.1
+ PRM + IRM	<b>58.9</b>	<b>25.4</b>
Llama2-7B: WizardMath-SFT	77.4	35.6
+ PRM	81.7	39.9
+ PRM + IRM	<b>84.1</b>	<b>43.5</b>
Mistral-7B: WizardMath-SFT	82.8	48.1
+ PRM	87.2	52.7
+ PRM + IRM	<b>90.7</b>	<b>55.4</b>

of our RLEIF method in enhancing mathematical reasoning capabilities across a range of problem difficulties, from grade to high school levels.

2) By employing diverse pre-trained models (i.e., GPT-2, Llama 2, Mistral, Qwen, DeepSeek) as base models, WizardMath demonstrated notable advancements on the GSM8k and MATH benchmarks. Specifically, WizardMath-Llama2-7B, based on Llama2-7B, improved performance by 69.5% on GSM8k and 41.0% on MATH. Similarly, WizardMath-GPT2-XL, built on GPT2-XL, achieved a 43.5% improvement on GSM8k and 18.5% on MATH, performing on par with Llama2-70B and outperforming GPT-3.5 on GSM8k. This demonstrates that our RLEIF method is equally effective for smaller models in enhancing mathematical reasoning capabilities, proving its scalability and robustness across various model backbones.

#### 4.4 ANALYSIS

##### The impact of training data size

We are curious about to how the training data size of different dataset construction methods impact the reasoning capacity of LLMs. Thus we conduct different number of training instances from ours evolved data and MetaMathQA to fine tune Mistral 7B. As shown in the Figure 2, Math Evol-Instruct achieves superior data efficiency. Specifically, our model constantly outperforms MetaMath by more than 3% ~ 6% on GSM8k and 15% ~ 20% on MATH under the same number of conditions. Our findings indicate that Math Evol-Instruct exhibits a higher potential upper bound compared to MetaMath, thus demonstrating the effectiveness of Evol-Instruct for math reasoning scenario.

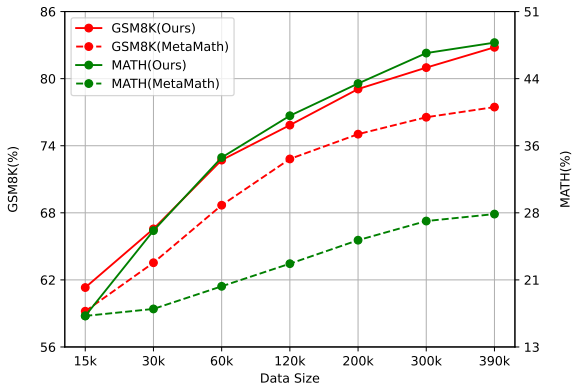


Figure 2: Accuracy of Mistral-7B fine-tuned in different sizes of augmentation data on GSM8K and MATH

##### The impact of PRM and IRM during PPO training

To verify the contributions of the instruction reward model and process-supervised reward model, we consider the following variants: (1) SFT + PRM: only use PRM in the PPO training. (2) SFT + PRM + IRM: use both IRM and PRM in the PPO training. As shown in Table 3, applying PRM alone for PPO training on GSM8k and MATH yields a 3%-4% improvement. When combined with IRM, an additional 2.5%-4% gain is observed. Thus, the integration of PRM and IRM results in a substantial overall improvement of 6%-8%. So, we can conclude that (1) PRM is crucial to WizardMath, since the variant with PRM significantly outperforms the SFT one without any PPO training (2) IRM also plays a key role in the success of reinforcement learning, as there is a remarkable improvement when

we combine PRM with IRM, further demonstrating the necessity of taking instruction’s quality into account and correcting false positives in the problem-solving process when we optimize the LLMs.

Table 4: The effect of different reward models during PPO training

Models	GSM8K	MATH
Llama2-7B: WizardMath-SFT	77.4	35.6
+ ORM (ours)	79.1	36.8
+ PRM800k	79.7	38.7
+ Math-Shepherd	80.3	38.2
+ PRM (ours)	<b>81.7</b>	<b>39.9</b>
Mistral-7B: WizardMath-SFT	82.8	48.1
+ ORM (ours)	84.6	49.6
+ PRM800k	85.4	50.8
+ Math-Shepherd	86.1	50.3
+ PRM (ours)	<b>87.2</b>	<b>52.7</b>

Table 5: Results of reinforcement learning combined with validation. The SFT and Reward models are trained based on Mistral-7B. The verifier is based on 256 sample outputs.

Generators	Verifiers	GSM8K	MATH
SFT	Self-Consistency	90.7	57.5
	ORM	93.0	58.3
	PRM	93.9	61.7
SFT + ORM	Self-Consistency	91.2	57.7
	ORM	93.4	59.4
	PRM	94.1	63.3
SFT + PRM	Self-Consistency	92.3	59.3
	ORM	94.1	60.8
	PRM	<b>95.2</b>	<b>64.7</b>

**The impact of Evol-Instruct turns.** Table 6 illustrates the impact of combining downward and upward evolution in SFT training. Two rounds of downward evolution improved GSM8k by 14.8% (74.5 vs. 59.7) and MATH by 19.6% (34.7 vs. 15.1) over the original. Three rounds of upward evolution yielded a 18.9% improvement on GSM8k (78.6 vs. 59.7) and a 27.4% improvement on MATH (42.5 vs. 15.1). Furthermore, combining downward evolution based on upward evolution resulted in an additional 2.6% improvement on GSM8k (81.2 vs. 78.6), a total improvement of 21.5% over the original. Similarly, a 1.9% improvement on MATH (46.5 vs. 42.5), a 31.4% total improvement. These results underscore the complementary and significant effectiveness of upward and downward evolution.

Table 6: Impact of different Downward and Upward Evol-Instruct turns on Mistral-7B SFT. *D-i* refers to the *i* round of downward evolution, whereas *U-i* denotes the *i* round of upward evolution. *Ori* is the original manually annotated 7.5k data of GSM8k and MATH.

Data	GSM8K							MATH						
	Ori	D-1	D-2	U-1	U-2	U-3	pass@1	Ori	D-1	D-2	U-1	U-2	U-3	pass@1
Ori	✓	x	x	x	x	x	59.7	✓	x	x	x	x	x	15.1
Math Evol	✓	✓	x	x	x	x	71.9	✓	✓	x	x	x	x	30.3
	✓	x	✓	x	x	x	70.5	✓	x	✓	x	x	x	28.7
	✓	x	x	✓	x	x	73.7	✓	x	x	✓	x	x	33.4
	✓	x	x	x	✓	x	71.6	✓	x	x	x	✓	x	32.6
	✓	x	x	x	x	✓	70.2	✓	x	x	x	x	✓	30.9
	✓	✓	✓	x	x	x	<b>74.5</b>	✓	✓	✓	x	x	x	<b>34.7</b>
	✓	x	x	✓	✓	x	77.1	✓	x	x	✓	✓	x	38.6
	✓	x	x	✓	✓	✓	<b>78.6</b>	✓	x	x	✓	✓	✓	<b>42.5</b>
	✓	✓	✓	✓	x	x	76.6	✓	✓	✓	✓	✓	✓	40.3
	✓	✓	✓	✓	✓	x	79.8	✓	✓	✓	✓	✓	x	44.6
✓	✓	✓	✓	✓	✓	<b>81.2</b>	✓	✓	✓	✓	✓	✓	<b>46.2</b>	

**ORM v.s. PRM; Human v.s. AI.** The Table 4 presents the performance of different answer reward methods for LLMs in terms of pass@1. As is shown: 1) Our step-by-step PRM significantly enhances the performance of both Llama and Mistral based SFT models. Specifically, the Mistral-7B powered by our PRM achieves 87.2% and 52.7% on GSM8k and MATH respectively. 2) PRM models consistently outperforms ORM on both GSM8k and MATH, indicating the effectiveness of step-by-step supervision. 3) The PRM trained on our fully AI-labeled data outperforms both the manually annotated PRM800k and Math-Shepherd, which utilizes MCTS tree search for annotation. When training WizardMath-Mistral-SFT with PPO, our PRM improves upon PRM800k by 1.8% and Math-Shepherd by 1.1% on GSM8k, while surpassing PRM800k by 1.9% and Math-Shepherd by 2.4% on MATH. This demonstrates powerful AI can also provide good process supervision quality, highlighting the effectiveness of utilizing AI to construct PRM training data.

**PRM as Verifier.** Table 5 presents the performance comparison of various generators with different verifiers on GSM8K and MATH in terms of pass@256. We find that: 1) PRM verifier consistently demonstrates superior performance compared to Self-Consistency and ORM. Specifically, our SFT + PRM generator, enhanced by the PRM verifier, achieves 95.2% and 64.7% accuracy on GSM8K and MATH respectively. 2) When compared to ORM, PRM exhibits a more significant advantage on the more challenging MATH dataset which aligns with the findings in (Uesato et al., 2022) and (Lightman et al., 2023). This can be attributed to the fact that GSM8K involves fewer and less complex steps in problem-solving than MATH. 3) Particularly, the generator with PRM PPO training



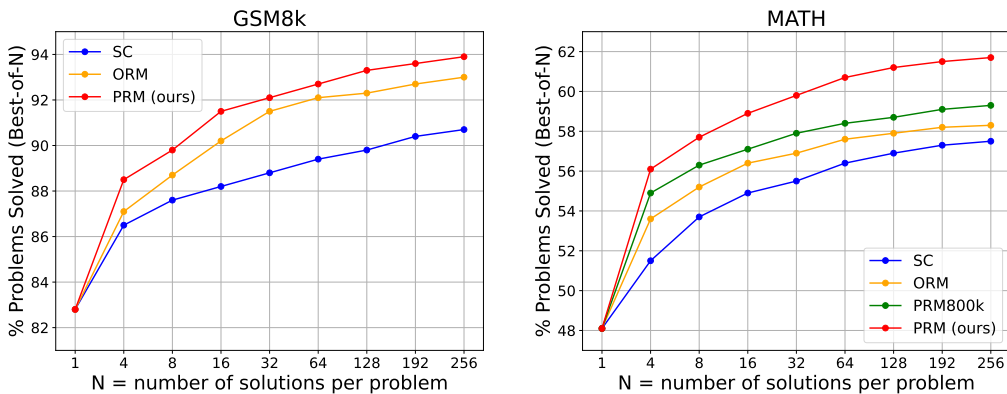


Figure 3: Performance of Mistral-7B SFT with different verification strategies.

surpasses those SFT and ORM PPO trained generators regardless of employing Self-Consistency, ORM, and the PRM verifiers. This further demonstrates the effectiveness of our PRM.

Figure 3 also shows the performance of different Verification strategies across a range of candidate numbers from 1 to 256 on two benchmarks. The main observations are as follows: 1) PRM verifiers consistently achieves superior performance compared to both ORM and majority voting, and this superiority becomes more evident as N increases. 2) For MATH benchmark, our PRM trained on the AI-annotated datasets slightly surpassed the human-annotated PRM800K.

**Performance of Out-of-Domain.** Table 7 presents the results of WizardMath on the 7 out-of-domain evaluation results covering K-12, college, and competition level math problems, highlighting the following salient observations: (1) With math Evol-Instruct and reinforcement learning, WizardMath consistently surpasses prior state-of-the-art open-source models (e.g. MetaMath, MathScale) across all scales, and achieves improvement of 5%-10% across 7 tasks on average. (2) The accuracy of WizardMath-Mistral is about 5.0% higher than WizardMath-Llama on the same size. Especially it exceeds GPT-3.5-Turbo (45.7 vs. 37.9) while being comparable to GPT-4. This also indicates that Mistral-7B has more potential in mathematical reasoning. (3) Especially on difficult benchmarks (i.e., College Math, AGIE Gaokao Math), WizardMath outperforms MetaMath by a significant margin. This demonstrates our model and RLEIF method has stronger robustness and better significant generalization ability for invisible mathematical problems.

Table 7: Performance of WizardMath on the 7 out-of-domain evaluation results covering K-12, college, and competition level math problems. The results of models in the table refer to MWPBENCH (Tang et al., 2024). “AGIE” stands for AGIEval. We report the models’ CoT pass@1 results on MwpBench without using any external python tool

Models	College Math	TAL	Math23k	Ape210k	Gaokao Bench Math	AGIE Gaokao Math	AGIE SAT Math	AVG
<i>Proprietary models</i>								
GPT-4	24.4	51.8	76.5	61.5	35.4	28.2	68.6	49.5
GPT-3.5-Turbo	21.6	42.9	62.5	44.0	23.2	15.3	55.8	37.9
<i>Models based on LLaMA-2 13B</i>								
LLaMA-2 13B	1.2	6.3	9.5	7.9	0.7	0.4	6.8	4.7
MAmmoTH-CoT	6.5	17.3	39.5	28.1	5.9	4.9	20.5	17.5
GAIR-Abel	7.9	21.1	42.2	27.8	7.0	4.9	30.3	20.2
MetaMath	10.1	25.4	48.6	31.6	9.6	5.6	38.2	24.2
MathScale 13B	20.4	38.1	61.1	43.7	20.0	12.3	55.8	35.9
WizardMath	22.9	43.3	70.3	50.8	33.1	25.7	64.7	44.4
<i>Models based on LLaMA-2 7B</i>								
LLaMA-2 7B	2.3	7.6	6.8	7.3	2.1	2.9	2.9	4.6
MAmmoTH-CoT	6.2	13.3	34.6	21.4	3.9	2.7	19.6	14.5
GAIR-Abel	6.6	18.3	35.4	24.5	4.3	4.4	23.5	16.7
MetaMath	9.4	22.5	44.0	29.9	5.9	5.1	36.2	21.9
MathScale 7B	20.9	35.2	59.0	41.8	19.6	12.6	57.8	35.3
WizardMath	21.2	40.2	67.3	46.1	28.9	18.7	62.7	40.7
<i>Models based on Mistral 7B</i>								
Mistral 7B	7.5	17.9	18.5	15.5	6.2	5.9	22.5	13.4
MetaMath Mistral	15.7	31.4	55.1	38.1	15.3	10.1	50.9	30.9
MathScale Mistral	21.8	39.9	64.4	46.0	21.4	14.3	57.8	37.9
WizardMath Mistral	24.8	44.8	71.2	52.6	37.2	24.5	64.7	45.7

**Employ Open-source Model to Math Evol-Instruct.** In Table 33, we investigate the use of open-source models (i.e., Llama-3-70B-Instruct) as a substitute for GPT-4 during the SFT stage for Evol Instruct, employing the same evolution strategy. The results demonstrate that WizardMath-

Table 9: A case study from GSM8k test set. We rate the response using PRM and ORM. Red text denotes the wrong reasoning steps which PRM successfully detected, but ORM failed.

486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503	<p><b>Question</b> <i>The farm has 30 cows and the zoo has 20 sheep. The zoo has twice as many cows as the farm does. The farm has half as many sheep as the zoo does. How many animals do the farm and zoo have combined? <b>Ground truth answer: 120</b></i></p> <hr/> <p><b>Steps</b></p> <p>Step 1: we know that the farm has 30 cows. PRM Score: 0.95.</p> <p>Step 2: we are told that the zoo has twice as many cows as the farm does. PRM Score: 0.97.</p> <p>Step 3: So, we can calculate the number of cows in the zoo by multiplying the number of cows in the farm by 2. So, <math>2 * 30 = 60</math>, So, the zoo has 60 cows. PRM Score: 0.97.</p> <p>Step 4: we know that the farm has half as many sheep as the zoo does. PRM Score: 0.96.</p> <p>Step 5: So, we can calculate the number of sheep in the farm by dividing the number of sheep in the zoo by 2. So, <math>20 / 2 = 10</math>. So, the farm has 10 sheep. PRM Score: 0.93.</p> <p>Step 6: Now, we can calculate the total number of animals in the farm and the zoo by adding the number of cows and the number of sheep in both places. PRM Score: 0.92.</p> <p>Step 7: So, <math>30</math> (farm's cows) + <math>60</math> (zoo's cows) + <math>10</math> (farm's sheep) = <math>100</math>. PRM Score: 0.11.</p> <p>Step 8: So, the farm and the zoo have a total of 100 animals combined. PRM Score: 0.10.</p> <p>Step 9: Therefore, the final answer is <span style="border: 1px solid red; padding: 2px;">100</span>. The answer is: 100. PRM Score: 0.06. ORM Score: 0.89.</p>
--	---

Llama3-Evol achieved a 33.8% improvement on GSM8k and a 30.6% improvement on MATH, indicating that the math evol instruct strategy remains effective on open-source models. However, compared to GPT-4 evolution, there is still a 5%-6% performance gap. Despite this, the strategy shows significant potential in balancing computational cost and accuracy.

Table 8: The impact of using open source models for Math-Evol and use Mistral-7B-v0.1 for SFT .

Models	GSM8k	MATH
Mistral-7B-v0.1	42.9	12.9
WizardMath-SFT-GPT-4-Evol	82.8	48.1
WizardMath-SFT-Llama3-Evol	76.7	43.5

#### 4.5 MORE DISCUSSION.

Due to limited space, we place more discussion in the appendix. (1.) Appendix A.4 explores the effect of math evol-instruct during the SFT and RL stages, showing that math evol-instruct is highly efficient in SFT and RL stages. (2.) Appendix A.5 explores the difference between Math Evol-Instruct and WizardLM Evol-Instruct, showing math evol-instruct is more efficient than WizardLM. (3.) Appendix A.8 explores the impact of the different round for upward and downward evol-instruct. (4.) Appendix A.9 explores the impact of the scoring aggregation strategy at each step of the PRM for RL training. (5.) Appendix A.11 explores the data contamination check to prevent data leakage.

#### 4.6 CASE STUDY

**Evol-Instruct.** The Examples 3 and 4 in the Appendix A.1 shows the prompt and corresponding cases of GSM8k and MATH instruction evolution, demonstrating that the evolved instructions exhibit more complexity and diversity than the original training set.

**PRM vs. ORM.** We present a comprehensive case study to illustrate the effectiveness of our PRM. As delineated in Table 9, PRM demonstrates precise performance on a challenge math problem from the GSM8k test set. Remarkably, our PRM effectively distinguished the incorrect solution, in the meanwhile the ORM struggled in this task. Furthermore, PRM demonstrated exceptional insight by accurately detecting the incorrect steps of the solution chosen by ORM, specifically the steps 7, 8, and 9. Subsequently, PRM also assigned lower score logits to these erroneous steps.

## 5 CONCLUSION

This paper introduces *WizardMath*, a mathematics model fine-tuned with *RLEIF*. The experimental results demonstrate that *WizardMath* achieves SOTA performance surpassing existing open-source LLMs on GSM8k and MATH from grade to high school problems. Notably, *WizardMath* 70B exhibits superior performance compared to some of the well-known proprietary LLMs, including ChatGPT-3.5, Claude Instant, PaLM-2, Gemini Pro. Furthermore, our preliminary exploration highlights the pivotal role of instruction evolution and process supervision in achieving exceptional performance.

## REFERENCES

- 540  
541  
542 Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. Large language models  
543 for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157*, 2024.
- 544 Avinash Anand, Mohit Gupta, Kritarth Prasad, Navya Singla, Sanjana Sanjeev, Jatin Kumar,  
545 Adarsh Raj Shivam, and Rajiv Ratn Shah. Mathify: Evaluating large language models on mathe-  
546 matical problem solving tasks. *arXiv preprint arXiv:2404.13099*, 2024.
- 547 Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,  
548 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv*  
549 *preprint arXiv:2305.10403*, 2023.
- 550  
551 Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q  
552 Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for  
553 mathematics. *arXiv preprint arXiv:2310.10631*, 2023.
- 554 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenhang Ge,  
555 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao  
556 Liu, Chengqiang Lu, K. Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi  
557 Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng  
558 Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Jian Yang, Shusheng Yang, Shusheng Yang, Bowen  
559 Yu, Yu Bowen, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xing Zhang, Yichang Zhang, Zhenru  
560 Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report.  
561 *ArXiv*, abs/2309.16609, 2023. URL [https://api.semanticscholar.org/CorpusID:  
562 263134555](https://api.semanticscholar.org/CorpusID:263134555).
- 563 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna  
564 Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness  
565 from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- 566  
567 DeepSeek-AI Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng,  
568 Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi  
569 Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He, Wen-Hui  
570 Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y. K. Li, Wenfeng Liang, Fangyun  
571 Lin, A. X. Liu, Bo Liu, Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu Lu, Shanghao Lu,  
572 Fuli Luo, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui Qu, Tongzheng Ren,  
573 Zehui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Jun-Mei Song, Xuecheng Su, Jingxiang  
574 Sun, Yaofeng Sun, Min Tang, Bing-Li Wang, Peiyi Wang, Shiyu Wang, Yaohui Wang, Yongji  
575 Wang, Tong Wu, Y. Wu, Xin Xie, Zhenda Xie, Ziwei Xie, Yi Xiong, Hanwei Xu, Ronald X Xu,  
576 Yanhong Xu, Dejian Yang, Yu mei You, Shuiping Yu, Xin yuan Yu, Bo Zhang, Haowei Zhang,  
577 Lecong Zhang, Liyue Zhang, Mingchuan Zhang, Minghu Zhang, Wentao Zhang, Yichao Zhang,  
578 Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shunfeng Zhou, Qihao Zhu, and Yuheng Zou.  
579 Deepseek llm: Scaling open-source language models with longtermism. *ArXiv*, abs/2401.02954,  
2024. URL <https://api.semanticscholar.org/CorpusID:266818336>.
- 580 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
581 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
582 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
583 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,  
584 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-  
585 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot  
586 learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan,  
587 and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual*  
588 *Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,*  
589 *2020, virtual*, 2020a. URL [https://proceedings.neurips.cc/paper/2020/hash/  
590 1457c0d6bfc4967418bfb8ac142f64a-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html).
- 591 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
592 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
593 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,  
Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott

- 594 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya  
595 Sutskever, and Dario Amodei. Language models are few-shot learners, 2020b.  
596
- 597 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,  
598 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:  
599 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 600 Jonathon Cai, Richard Shin, and Dawn Song. Making neural programming architectures generalize  
601 via recursion. *arXiv preprint arXiv:1704.06611*, 2017.  
602
- 603 Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. Scaling synthetic data creation with  
604 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*, 2024.  
605
- 606 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
607 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,  
608 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
609 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,  
610 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios  
611 Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino,  
612 Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,  
613 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa,  
614 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob  
615 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating  
large language models trained on code, 2021.
- 616 Zhaorun Chen, Zhuokai Zhao, Zhihong Zhu, Ruiqi Zhang, Xiang Li, Bhiksha Raj, and Huaxiu Yao.  
617 Autoprnm: Automating procedural supervision for multi-step reasoning via controllable question  
618 decomposition. *arXiv preprint arXiv:2402.11452*, 2024a.  
619
- 620 Zhipeng Chen, Kun Zhou, Wayne Xin Zhao, Junchen Wan, Fuzheng Zhang, Di Zhang, and Ji-Rong  
621 Wen. Improving large language models via fine-grained reinforcement learning with minimum  
622 editing constraint. *arXiv preprint arXiv:2401.06081*, 2024b.
- 623 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning  
624 converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*,  
625 2024c.  
626
- 627 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
628 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An  
629 open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://vicuna.lmsys.org>.  
630
- 631 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
632 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,  
633 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam  
634 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James  
635 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Lev-  
636 skaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin  
637 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph,  
638 Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M.  
639 Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon  
640 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark  
641 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean,  
642 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways, 2022.
- 643 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
644 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve  
645 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.  
646
- 647 Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large  
language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.

- 648 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning  
649 of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- 650
- 651 Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and  
652 Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv*  
653 *preprint arXiv:2304.06767*, 2023.
- 654 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
655 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
656 *arXiv preprint arXiv:2407.21783*, 2024.
- 657
- 658 Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun  
659 Wang. Alphazero-like tree-search can guide large language model decoding and training. *arXiv*  
660 *preprint arXiv:2309.17179*, 2023.
- 661 Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz,  
662 Philipp Christian Petersen, Alexis Chevalier, and Julius Berner. Mathematical capabilities of  
663 chatgpt. *arXiv preprint arXiv:2301.13867*, 2023.
- 664 Jiayi Fu, Lei Lin, Xiaoyang Gao, Pengli Liu, Zhengzong Chen, Zhirui Yang, Shengnan Zhang, Xue  
665 Zheng, Yan Li, Yuliang Liu, et al. Kwaiyimath: Technical report. *arXiv preprint arXiv:2310.07488*,  
666 2023a.
- 667
- 668 Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting  
669 for multi-step reasoning. *arXiv preprint arXiv:2210.00720*, 2022.
- 670 Yao Fu, Litu Ou, Mingyu Chen, Yuhao Wan, Hao Peng, and Tushar Khot. Chain-of-thought hub:  
671 A continuous effort to measure large language models’ reasoning performance. *arXiv preprint*  
672 *arXiv:2305.17306*, 2023b.
- 673
- 674 Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle  
675 use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions*  
676 *of the Association for Computational Linguistics*, 9:346–361, 2021. doi: 10.1162/tacl\_a\_00370.  
677 URL <https://aclanthology.org/2021.tacl-1.21>.
- 678 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu  
679 Feng, Hanlin Zhao, Hanyu Lai, et al. Chatglm: A family of large language models from glm-130b  
680 to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.
- 681 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and  
682 Weizhu Chen. Tora: A tool-integrated reasoning agent for mathematical problem solving, 2023.
- 683
- 684 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint*  
685 *arXiv:1410.5401*, 2014.
- 686
- 687 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
688 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*  
689 *preprint arXiv:2103.03874*, 2021.
- 690 Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to  
691 solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 Conference*  
692 *on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 523–533, Doha, Qatar,  
693 October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1058. URL  
694 <https://aclanthology.org/D14-1058>.
- 695 Xuhan Huang, Qingning Shen, Yan Hu, Anningzhe Gao, and Benyou Wang. Mamo: a mathematical  
696 modeling benchmark with solvers. *arXiv preprint arXiv:2405.13144*, 2024a.
- 697
- 698 Yiming Huang, Xiao Liu, Yeyun Gong, Zhibin Gou, Yelong Shen, Nan Duan, and Weizhu Chen.  
699 Key-point-driven data synthesis with its enhancement on mathematical reasoning. *arXiv preprint*  
700 *arXiv:2403.02333*, 2024b.
- 701 Shima Imani, Liang Du, and Harsh Shrivastava. Mathprompter: Mathematical reasoning using large  
language models. *arXiv preprint arXiv:2303.05398*, 2023.

- 702 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
703 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
704 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.  
705
- 706 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
707 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand,  
708 Gianna Lengyel, Guillaume Bour, Guillaume Lample, L elio Renard Lavaud, Lucile Saulnier, Marie-  
709 Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le  
710 Scao, Th eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth ee Lacroix, and William El Sayed.  
711 Mixtral of experts, 2024a.
- 712 Weisen Jiang, Han Shi, Longhui Yu, Zhengying Liu, Yu Zhang, Zhenguo Li, and James Kwok.  
713 Forward-backward reasoning in large language models for mathematical verification. In *Findings*  
714 *of the Association for Computational Linguistics ACL 2024*, pp. 6647–6661, 2024b.
- 715 Zhanming Jie, Jierui Li, and Wei Lu. Learning to reason deductively: Math word problem solving as  
716 complex relation extraction. *arXiv preprint arXiv:2203.10316*, 2022.  
717
- 718 Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang.  
719 Parsing algebraic word problems into equations. *Transactions of the Association for Computational*  
720 *Linguistics*, 3:585–597, 2015.  
721
- 722 Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. MAWPS:  
723 A math word problem repository. In *Proceedings of the 2016 Conference of the North American*  
724 *Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp.  
725 1152–1157, San Diego, California, June 2016. Association for Computational Linguistics. doi:  
726 10.18653/v1/N16-1136. URL <https://aclanthology.org/N16-1136>.
- 727 Yihuai Lan, Lei Wang, Qiyuan Zhang, Yunshi Lan, Bing Tian Dai, Yan Wang, Dongxiang Zhang,  
728 and Ee-Peng Lim. Mwptoolkit: an open-source framework for deep learning-based math word  
729 problem solvers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp.  
730 13188–13190, 2022.  
731
- 732 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-  
733 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative  
734 reasoning problems with language models. *arXiv preprint arXiv:2206.14858*, 2022.
- 735 Chen Li, Weiqi Wang, Jingcheng Hu, Yixuan Wei, Nanning Zheng, Han Hu, Zheng Zhang, and  
736 Houwen Peng. Common 7b language models already possess strong math capabilities. *arXiv*  
737 *preprint arXiv:2403.04706*, 2024a.  
738
- 739 Chengpeng Li, Zheng Yuan, Guanting Dong, Keming Lu, Jiancan Wu, Chuanqi Tan, Xiang Wang,  
740 and Chang Zhou. Query and response augmentation cannot help out-of-domain math reasoning  
741 generalization. *ArXiv*, abs/2310.05506, 2023a. URL <https://api.semanticscholar.org/CorpusID:263830207>.  
742
- 743 Chengpeng Li, Guanting Dong, Mingfeng Xue, Ru Peng, Xiang Wang, and Dayiheng Liu. Dotamath:  
744 Decomposition of thought with code assistance and self-correction for mathematical reasoning.  
745 *arXiv preprint arXiv:2407.04078*, 2024b.  
746
- 747 Chengtao Li, Daniel Tarlow, Alexander L. Gaunt, Marc Brockschmidt, and Nate Kushman. Neural  
748 program lattices. In *International Conference on Learning Representations*, 2016. URL <https://api.semanticscholar.org/CorpusID:34816748>.  
749
- 750 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif  
751 Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in  
752 ai4maths with 860k pairs of competition math problems and solutions. 2024c.  
753
- 754 Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. Gsm-plus: A comprehensive  
755 benchmark for evaluating the robustness of llms as mathematical problem solvers. *arXiv preprint*  
*arXiv:2402.19255*, 2024d.

- 756 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,  
757 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with  
758 you! *arXiv preprint arXiv:2305.06161*, 2023b.
- 759 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, B. Chen, Jian-Guang Lou, and Weizhu Chen. Making  
760 language models better reasoners with step-aware verifier. In *Annual Meeting of the Association  
761 for Computational Linguistics*, 2022a. URL [https://api.semanticscholar.org/  
762 CorpusID:259370847](https://api.semanticscholar.org/CorpusID:259370847).
- 764 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. On the  
765 advance of making language models better reasoners. *arXiv preprint arXiv:2206.02336*, 2022b.
- 766 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making  
767 language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual  
768 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5315–5333,  
769 Toronto, Canada, July 2023c. Association for Computational Linguistics. doi: 10.18653/v1/2023.  
770 acl-long.291. URL <https://aclanthology.org/2023.acl-long.291>.
- 772 Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards  
773 general text embeddings with multi-stage contrastive learning. *ArXiv*, abs/2308.03281, 2023d.  
774 URL <https://api.semanticscholar.org/CorpusID:260682258>.
- 775 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan  
776 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint  
777 arXiv:2305.20050*, 2023.
- 778 Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Yelong Shen, Ruochen Xu, Chen Lin, Yujiu  
779 Yang, Jian Jiao, Nan Duan, et al. Rho-1: Not all tokens are what you need. *arXiv preprint  
780 arXiv:2404.07965*, 2024.
- 782 Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale genera-  
783 tion: Learning to solve and explain algebraic word problems. *ACL*, 2017.
- 784 Bingbin Liu, Sebastien Bubeck, Ronen Eldan, Janardhan Kulkarni, Yuanzhi Li, Anh Nguyen, Rachel  
785 Ward, and Yi Zhang. Tinygsm: achieving > 80% on gsm8k with small language models. *arXiv  
786 preprint arXiv:2312.09241*, 2023.
- 788 Haoxiong Liu and Andrew Chi-Chih Yao. Augmenting math word problems via iterative question  
789 composing. *arXiv preprint arXiv:2401.09003*, 2024.
- 790 Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang.  
791 On llms-driven synthetic data generation, curation, and evaluation: A survey. *arXiv preprint  
792 arXiv:2406.15126*, 2024.
- 793 Pan Lu, Liang Qiu, Wenhao Yu, Sean Welleck, and Kai-Wei Chang. A survey of deep learning for  
794 mathematical reasoning. *arXiv preprint arXiv:2212.10535*, 2022.
- 796 Zimu Lu, Aojun Zhou, Houxing Ren, Ke Wang, Weikang Shi, Junting Pan, Mingjie Zhan, and  
797 Hongsheng Li. Mathgenie: Generating synthetic data with question back-translation for enhancing  
798 mathematical reasoning of llms. *arXiv preprint arXiv:2402.16352*, 2024.
- 799 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing  
800 Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with  
801 evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
- 803 Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality  
804 in abstractive summarization. *arXiv preprint arXiv:2005.00661*, 2020.
- 805 Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-  
806 free reward. *arXiv preprint arXiv:2405.14734*, 2024.
- 807  
808 Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Codas,  
809 Yadong Lu, Wei-ge Chen, Olga Vrousgos, Corby Rosset, et al. Agentinstruct: Toward generative  
teaching with agentic flows. *arXiv preprint arXiv:2407.03502*, 2024a.

- 810 Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the  
811 potential of slms in grade school math. *arXiv preprint arXiv:2402.14830*, 2024b.
- 812
- 813 Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed  
814 Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint*  
815 *arXiv:2306.02707*, 2023.
- 816 Reiiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher  
817 Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted  
818 question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- 819
- 820 Ansong Ni, Jeevana Priya Inala, Chenglong Wang, Alex Polozov, Christopher Meek, Dragomir  
821 Radev, and Jianfeng Gao. Learning math reasoning from self-sampled correct and partially-correct  
822 solutions. In *The Eleventh International Conference on Learning Representations*, 2022.
- 823 Xinzhe Ni, Yeyun Gong, Zhibin Gou, Yelong Shen, Yujia Yang, Nan Duan, and Weizhu Chen. Explor-  
824 ing the mystery of influential data for mathematical reasoning. *arXiv preprint arXiv:2404.01067*,  
825 2024.
- 826
- 827 Eric Nichols, Leo Gao, and Randy Gomez. Collaborative storytelling with large-scale neural language  
828 models. In *Proceedings of the 13th ACM SIGGRAPH Conference on Motion, Interaction and*  
829 *Games*, pp. 1–10, 2020.
- 830 OpenAI. Gpt-4 technical report, 2023.
- 831
- 832 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,  
833 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser  
834 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan  
835 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In  
836 *NeurIPS*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/  
837 b1efde53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html).
- 838 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math  
839 word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the*  
840 *Association for Computational Linguistics: Human Language Technologies*, pp. 2080–2094, 2021.
- 841 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli,  
842 Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb  
843 dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv*  
844 *preprint arXiv:2306.01116*, 2023.
- 845
- 846 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
847 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*  
848 *in Neural Information Processing Systems*, 36, 2024.
- 849 Scott Reed and Nando De Freitas. Neural programmer-interpreters. *arXiv preprint arXiv:1511.06279*,  
850 2015.
- 851
- 852 Subhro Roy and Dan Roth. Solving general arithmetic word problems. In *Proceedings of the*  
853 *2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1743–1752, Lisbon,  
854 Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1202.  
855 URL <https://aclanthology.org/D15-1202>.
- 856 Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman  
857 Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-  
858 parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- 859
- 860 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li, Yu Wu,  
861 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language  
862 models. *arXiv preprint arXiv:2402.03300*, 2024.
- 863 Jianhao Shen, Yichun Yin, Lin Li, Lifeng Shang, Xin Jiang, Ming Zhang, and Qun Liu. Generate &  
rank: A multi-task framework for math word problems. *arXiv preprint arXiv:2109.03034*, 2021.



- 864 Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang.  
865 Preference ranking optimization for human alignment. *arXiv preprint arXiv:2306.17492*, 2023.  
866
- 867 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,  
868 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in*  
869 *Neural Information Processing Systems*, 33:3008–3021, 2020.
- 870 Zhaochen Su, Jun Zhang, Tong Zhu, Xiaoye Qu, Juntao Li, Min Zhang, and Yu Cheng. Timo:  
871 Towards better temporal reasoning for language models. *arXiv preprint arXiv:2406.14192*, 2024.  
872
- 873 Zhiqing Sun, Longhui Yu, Yikang Shen, Weiyang Liu, Yiming Yang, Sean Welleck, and Chuang  
874 Gan. Easy-to-hard generalization: Scalable alignment beyond human supervision. *arXiv preprint*  
875 *arXiv:2403.09472*, 2024.
- 876 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question  
877 answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of*  
878 *the North American Chapter of the Association for Computational Linguistics: Human Language*  
879 *Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June  
880 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421>.  
881
- 882 Zhengyang Tang, Xingxing Zhang, Benyou Wan, and Furu Wei. Mathscales: Scaling instruction  
883 tuning for mathematical reasoning. *arXiv preprint arXiv:2403.02884*, 2024.  
884
- 885 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
886 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
887 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- 888 Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia,  
889 Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science.  
890 *arXiv preprint arXiv:2211.09085*, 2022.
- 891 Gemini Team. Gemini: A family of highly capable multimodal models, 2023.  
892
- 893 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu  
894 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable  
895 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 896 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,  
897 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models  
898 based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.  
899
- 900 Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. Toward self-  
901 improvement of llms via imagination, searching, and criticizing. *arXiv preprint arXiv:2404.12253*,  
902 2024.
- 903 Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu, and Junxian He. Dart-math: Difficulty-aware  
904 rejection tuning for mathematical problem-solving. *arXiv preprint arXiv:2407.13690*, 2024.  
905
- 906 Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Git-  
907 man. Openmathinstruct-1: A 1.8 million math instruction tuning dataset. *arXiv preprint*  
908 *arXiv:2402.10176*, 2024.
- 909 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
910 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
911 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 912 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
913 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
914 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.  
915
- 916 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia  
917 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and  
outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.

- 918 Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song,  
919 Mingjie Zhan, and Hongsheng Li. Mathcoder: Seamless code integration in llms for enhanced  
920 mathematical reasoning. *arXiv preprint arXiv:2310.03731*, 2023a.
- 921  
922 Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim.  
923 Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language  
924 models. *arXiv preprint arXiv:2305.04091*, 2023b.
- 925 Peiyi Wang, Lei Li, Liang Chen, Feifan Song, Binghuai Lin, Yunbo Cao, Tianyu Liu, and Zhifang Sui.  
926 Making large language models better reasoners with alignment. *ArXiv*, abs/2309.02144, 2023c.  
927 URL <https://api.semanticscholar.org/CorpusID:261558535>.
- 928 Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and  
929 Zhifang Sui. Large language models are not fair evaluators. *ArXiv*, abs/2305.17926, 2023d. URL  
930 <https://api.semanticscholar.org/CorpusID:258960339>.
- 931  
932 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang  
933 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In  
934 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume*  
935 *1: Long Papers)*, pp. 9426–9439, 2024a.
- 936 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-  
937 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.  
938 *arXiv preprint arXiv:2203.11171*, 2022.
- 939 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha  
940 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language  
941 models. In *ICLR 2023*, 2023e. URL <https://arxiv.org/abs/2203.11171>.
- 942  
943 Yue Wang, Weishi Wang, Shafiq R. Joty, and Steven C. H. Hoi. Codet5: Identifier-aware unified pre-  
944 trained encoder-decoder models for code understanding and generation. In Marie-Francine Moens,  
945 Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference*  
946 *on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta*  
947 *Cana, Dominican Republic, 7-11 November, 2021*, pp. 8696–8708. Association for Computational  
948 Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.685. URL [https://doi.org/10.](https://doi.org/10.18653/v1/2021.emnlp-main.685)  
949 [18653/v1/2021.emnlp-main.685](https://doi.org/10.18653/v1/2021.emnlp-main.685).
- 950 Zengzhi Wang, Rui Xia, and Pengfei Liu. Generative ai for math: Part i–mathpile: A billion-token-  
951 scale pretraining corpus for math. *arXiv preprint arXiv:2312.17120*, 2023f.
- 952  
953 Zihan Wang, Yunxuan Li, Yuexin Wu, Liangchen Luo, Le Hou, Hongkun Yu, and Jingbo Shang.  
954 Multi-step problem solving through a verifier: An empirical analysis on model-induced process  
955 supervision. *arXiv preprint arXiv:2402.02658*, 2024b.
- 956 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny  
957 Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint*  
958 *arXiv:2201.11903*, 2022.
- 959 Shaohua Wu, Xudong Zhao, Shenling Wang, Jiangang Luo, Lingjun Li, Xi Chen, Bing Zhao,  
960 Wei Wang, Tong Yu, Rongguo Zhang, et al. Yuan 2.0: A large language model with localized  
961 filtering-based attention. *arXiv preprint arXiv:2311.15786*, 2023a.
- 962  
963 Yanan Wu, Jie Liu, Xingyuan Bu, Jiaheng Liu, Zhanhui Zhou, Yuanxing Zhang, Chenchen Zhang,  
964 Zhiqi Bai, Haibin Chen, Tiezheng Ge, et al. Conceptmath: A bilingual concept-wise benchmark  
965 for measuring mathematical reasoning of large language models. *arXiv preprint arXiv:2402.14660*,  
966 2024.
- 967 Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A.  
968 Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better  
969 rewards for language model training. *ArXiv*, abs/2306.01693, 2023b. URL [https://api.](https://api.semanticscholar.org/CorpusID:259064099)  
970 [semanticscholar.org/CorpusID:259064099](https://api.semanticscholar.org/CorpusID:259064099).
- 971 Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. Evaluating mathematical  
reasoning beyond accuracy. *arXiv preprint arXiv:2404.05692*, 2024.

- 972 Wei Xiong, Chengshuai Shi, Jiaming Shen, Aviv Rosenberg, Zhen Qin, Daniele Calandriello, Misha  
973 Khalman, Rishabh Joshi, Bilal Piot, Mohammad Saleh, et al. Building math agents with multi-turn  
974 iterative preference learning. *arXiv preprint arXiv:2409.02392*, 2024.
- 975  
976 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin  
977 Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv  
978 preprint arXiv:2304.12244*, 2023.
- 979 Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan,  
980 Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint  
981 arXiv:2309.10305*, 2023.
- 982  
983 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-  
984 hong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical  
985 expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- 986 Fei Yu, Anningzhe Gao, and Benyou Wang. Outcome-supervised verifiers for planning in mathe-  
987 matical reasoning. *ArXiv*, abs/2311.09724, 2023a. URL <https://api.semanticscholar.org/CorpusID:265221057>.
- 988  
989 Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language models are super mario:  
990 Absorbing abilities from homologous models as a free lunch. In *Forty-first International Conference  
991 on Machine Learning*, 2024.
- 992  
993 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok,  
994 Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical  
995 questions for large language models, 2023b.
- 996  
997 Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. Scaling  
998 relationship on learning mathematical reasoning with large language models. *arXiv preprint  
999 arXiv:2308.01825*, 2023a.
- 1000  
1001 Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, and Songfang Huang. How well do large  
1002 language models perform in arithmetic tasks? *arXiv preprint arXiv:2304.02015*, 2023b.
- 1003  
1004 Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf:  
1005 Rank responses to align language models with human feedback without tears. *arXiv preprint  
1006 arXiv:2304.05302*, 2023c.
- 1007  
1008 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.  
1009 Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint  
1010 arXiv:2309.05653*, 2023.
- 1011  
1012 Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen. Mammoth2: Scaling instructions from the  
1013 web. *arXiv preprint arXiv:2405.03548*, 2024.
- 1014  
1015 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with  
1016 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- 1017  
1018 Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,  
1019 Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint  
1020 arXiv:2210.02414*, 2022.
- 1021  
1022 Liang Zeng, Liangjun Zhong, Liang Zhao, Tianwen Wei, Liu Yang, Jujie He, Cheng Cheng, Rui Hu,  
1023 Yang Liu, Shuicheng Yan, et al. Skywork-math: Data scaling laws for mathematical reasoning in  
1024 large language models—the story goes on. *arXiv preprint arXiv:2407.08348*, 2024.
- 1025  
1026 Dan Zhang, Sining Zhou, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts\*: Llm self-training  
1027 via process reward guided tree search. *arXiv preprint arXiv:2406.03816*, 2024a.
- 1028  
1029 Di Zhang, Jiatong Li, Xiaoshui Huang, Dongzhan Zhou, Yuqiang Li, and Wanli Ouyang. Accessing  
1030 gpt-4 level mathematical olympiad solutions via monte carlo tree self-refine with llama-3 8b. *arXiv  
1031 preprint arXiv:2406.07394*, 2024b.

- 1026 Jiaxin Zhang, Zhongzhi Li, Mingliang Zhang, Fei Yin, Chenglin Liu, and Yashar Moshfeghi. Geoeval:  
1027 benchmark for evaluating llms and multi-modal models on geometry problem-solving. *arXiv*  
1028 *preprint arXiv:2402.10104*, 2024c.
- 1029 Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming  
1030 Lu. Self-contrast: Better reflection through inconsistent solving perspectives. *arXiv preprint*  
1031 *arXiv:2401.02009*, 2024d.
- 1032
- 1033 Yifan Zhang, Yifan Luo, Yang Yuan, and Andrew C Yao. Autonomous data selection with language  
1034 models for mathematical texts. In *ICLR 2024 Workshop on Navigating and Addressing Data*  
1035 *Problems for Foundation Models*, 2024e.
- 1036
- 1037 Xu Zhao, Yuxi Xie, Kenji Kawaguchi, Junxian He, and Qizhe Xie. Automatic model selection with  
1038 large language models for reasoning. *arXiv preprint arXiv:2305.14333*, 2023.
- 1039
- 1040 Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo Li, and Yu Li. Progressive-hint prompting  
1041 improves reasoning in large language models. *arXiv preprint arXiv:2304.09797*, 2023.
- 1042
- 1043 Denny Zhou, Nathanael Scharli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schu-  
1044 rrmans, Olivier Bousquet, Quoc Le, and Ed Huai hsin Chi. Least-to-most prompting en-  
1045 ables complex reasoning in large language models. *ArXiv*, abs/2205.10625, 2022. URL  
<https://api.semanticscholar.org/CorpusID:248986239>.
- 1046
- 1047 Kun Zhou, Beichen Zhang, Jiapeng Wang, Zhipeng Chen, Wayne Xin Zhao, Jing Sha, Zhichao Sheng,  
1048 Shijin Wang, and Ji-Rong Wen. Jiuzhang3. 0: Efficiently improving mathematical reasoning by  
training small data synthesis models. *arXiv preprint arXiv:2405.14365*, 2024.
- 1049
- 1050 Xinyu Zhu, Junjie Wang, Lin Zhang, Yuxiang Zhang, Yongfeng Huang, Ruyi Gan, Jiaying Zhang, and  
1051 Yujiu Yang. Solving math word problems via cooperative reasoning induced language models. In  
1052 *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume*  
1053 *1: Long Papers)*. Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.acl-long.  
245. URL <https://doi.org/10.18653/v1/2023.acl-long.245>.
- 1054
- 1055 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul  
1056 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*  
1057 *preprint arXiv:1909.08593*, 2019.
- 1058
- 1059
- 1060
- 1061
- 1062
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079

## A APPENDIX

### A.1 MATH EVOLUTION PROMPTS

#### Example 1: Upward Evolution Prompt

**Step 1:** Understand the core concept and structure of the "#Instruction#". Identify the key elements such as variables, conditions, participants, actions, or processes that can be manipulated to increase complexity. Also, recognize the theme of the instruction and ensure it remains consistent throughout the evolution.

**Step 2:** Formulate a comprehensive plan to increment the complexity of the "#Instruction#" based on the identified elements in Step 1. The plan should involve modifying or expanding at least three components from the list. It is crucial to ensure that all components in the instruction are logically interconnected and that the complexity increase is coherent and justified. The plan should avoid introducing variables or conditions without clear criteria for determining their values or without contributing to the overall complexity. In this step, consider adding more real-world constraints and dependencies between variables to make the problem more challenging. And you can also add more constraints, concretizing, increasing reasoning.

**Step 3:** Implement the plan step by step to create the "#Rewritten Instruction#". Ensure the rewritten instruction maintains a logical sequence and avoids ambiguity or confusion. If additional variables or conditions are introduced, provide clear and unambiguous methods or criteria for determining their values. The "#Rewritten Instruction#" should not exceed the original "#Instruction#" by more than 30 words to ensure readability and comprehension.

**Step 4:** Review the "#Rewritten Instruction#" thoroughly to identify any unreasonable elements or inconsistencies. Make sure the "#Rewritten Instruction#" is a more complex version of the "#Instruction#", and that it accurately reflects the intended increase in complexity. Adjust any part of the instruction that may lead to misunderstanding or ambiguity, and provide the "#Finally Rewritten Instruction#" without any supplementary explanation.

Please reply strictly in the following format:

Step 1

#Elements Identified#:

Step 2

#Plan#:

Step 3

#Rewritten Instruction#:

Step 4

#Finally Rewritten Instruction#:

#Instruction#:

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

### Example 2: Downward Evolution Prompt

**Step 1:** Understand the "#Instruction#" and identify all the components that can be modified to decrease complexity, so that it makes the instruction easier. These components can be variables, conditions, participants, actions, etc. The key is to keep the core scenario unchanged while ensuring that any new elements introduced do not cause ambiguity or confusion.

**Step 2:** Develop a comprehensive plan to decrease the complexity of the "#Instruction#" based on the components identified in Step 1. The plan should involve modifying at least three components from the list. It is important to ensure that all components in the instruction are logically interconnected and that the complexity decrease is justifiable. The plan should avoid introducing variables or conditions without clear criteria for determining their values. Our goal is revising high difficulty questions to lower difficulty, or producing a new and easier question with another different topic.

**Step 3:** Implement the plan step by step to create the "#Rewritten Instruction#". Make sure the rewritten instruction maintains a logical sequence and avoids ambiguity or confusion. If additional variables or conditions are introduced, provide clear and unambiguous methods or criteria for determining their values. The "#Rewritten Instruction#" should not exceed the original "#Instruction#" by more than 20 words.

**Step 4:** Review the "#Rewritten Instruction#" thoroughly to identify any unreasonable elements. Make sure the "#Rewritten Instruction#" is a easier version of the "#Instruction#". Adjust any part of the instruction that may lead to misunderstanding or ambiguity, and provide the "#Finally Rewritten Instruction#" without any explanation.

Please reply strictly in the following format:

Step 1

#Elements Identified#:

Step 2

#Plan#:

Step 3

#Rewritten Instruction#:

Step 4

#Finally Rewritten Instruction#:

#Instruction#:

**Example 3: GSM8k Evol Instruction Case**

**Original Instruction 1:** Bill is trying to decide whether to make blueberry muffins or raspberry muffins. Blueberries cost \$5.00 per 6 ounce carton and raspberries cost \$3.00 per 8 ounce carton. If Bill is going to make 4 batches of muffins, and each batch takes 12 ounces of fruit, how much money would he save by using raspberries instead of blueberries?

**Evol Instruction 1:** Bill and Jane are contemplating between blueberry and raspberry muffins. Blueberries are \$5.00 for a 6 ounce carton, with a 20% bulk discount. Raspberries are \$3.00 for an 8 ounce carton. If they each make 6 batches of muffins, with each batch requiring 12 ounces of fruit, calculate the total money they would save by choosing raspberries over the discounted blueberries, given Jane's inclination towards raspberries.

**Original Instruction 2:** A snake's head is one-tenth its length. If a snake is 10 feet long, calculate the length of the rest of its body minus the head.

**Evol Instruction 2:** Given a snake's head is a certain fraction of its total length, and the snake's total length is a positive integer, determine the length of the snake's head by multiplying the total length by the fraction. Subtract this value from the total length to calculate the length of the rest of the snake's body.

**Original Instruction 3:** Thomas is training at the gym to prepare for a competition. He trained for 5 hours every day for a month (30 days). If he continues to train for the next 12 days, how many hours will he spend on training in total?

**Evol Instruction 3:** Thomas and James are preparing for a competition by training at the gym. They trained for 5 hours daily for a month (30 days), excluding a rest day each week. If they persist in training for the subsequent 12 days, adding an extra hour of training each week, what will be the total hours they have spent training?

**Original Instruction 4:** Travis is hired to take 638 bowls from the factory to the home goods store. The home goods store will pay the moving company a \$100 fee, plus \$3 for every bowl that is delivered safely. Travis must pay the home goods store \$4 each for any bowls that are lost or broken. If 12 bowls are lost, 15 bowls are broken, and the rest are delivered safely, how much should Travis be paid?

**Evol Instruction 4:** Travis and his team are tasked with moving 1000 bowls and 500 plates from the factory to a home goods store. The store agrees to pay a \$200 fee, plus \$4 for each safely delivered bowl and \$2 for each plate. However, Travis must compensate the store \$5 for each lost or broken bowl and \$3 for each plate. If they lose 20 bowls and 10 plates, and break 25 bowls and 15 plates, how much should the store pay Travis and his team?

**Original Instruction 5:** Gary is buying chlorine for his rectangular pool, which is 10 feet long, 8 feet wide, and 6 feet deep. Gary needs to buy one quart of chlorine for every 120 cubic feet of water in his pool. If chlorine costs \$3 a quart, how much does Gary spend on chlorine?

**Evol Instruction 5:** Gary and John are purchasing chlorine for their cylindrical pools, with diameters of 12 feet and 10 feet, and depths of 8 feet and 6 feet respectively. They require one quart of chlorine per 100 cubic feet of pool water. Given that chlorine is priced at \$4 per quart, calculate the total expenditure on chlorine for both Gary and John.

**Original Instruction 6:** Ken likes to bike when it's raining and can cycle 30 miles in 20 minutes during this time. However, when it's snowing Ken can't stand the cold and can only cycle 10 miles in 20 minutes. If it rains 3 times and snows 4 times in one week, how many miles did Ken reach if he cycles 1 hour a day?

**Evol Instruction 6:** In varying weather conditions, Ken's biking speed differs. He can cycle 30 miles in 20 minutes when it's raining, 10 miles in 20 minutes when it's snowing, and 20 miles in 20 minutes on sunny days. In a week, if it rains 4 times, snows 3 times, and is sunny 2 times, and Ken cycles for 1.5 hours each day, how many miles did he cover? Remember, after cycling for an hour, his speed decreases by 10%.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

#### Example 4: MATH Evol Instruction Case

**Original Instruction 1:** Find the smallest positive integer whose cube ends in 888.

**Evol Instruction 1:** Determine the least positive whole number, denoted by 'x', whose cube terminates in 888 and is divisible by 3. Verify the result by checking the divisibility of the cube by 9.

**Original Instruction 2:** The sum of all the positive factors of integer  $x$  is 24. If one of the factors is 3, what is the value of  $x$ ?

**Evol Instruction 2:** Given that the summation of all positive factors of an integer  $x$  is 24, and considering  $x$  is a positive integer divisible by 3 with one of its factors being 3, determine the value of  $x$  by first calculating the variable  $S$  representing the sum of factors, and then solving for  $x$ .

**Original Instruction 3:** What is  $2^{-1} + 2^{-2} + 2^{-3} + 2^{-4} + 2^{-5} + 2^{-6} \pmod{13}$ ? Express your answer as an integer from 0 to 12, inclusive.

**Evol Instruction 3:** Let  $S$  be the sum of the series  $2^{-1} + 2^{-2} + 2^{-3} + 2^{-4} + 2^{-5} + 2^{-6}$ . Calculate  $S$  by finding the sum of each term, then determine the value of  $S \pmod{13}$ . Utilize the properties of modular arithmetic and provide a step-by-step solution. Express the final answer as an integer from 0 to 12, inclusive.

**Original Instruction 4:** Find the greatest common divisor of 40304 and 30203.

**Evol Instruction 4:** Determine the greatest common divisor of the integers 40304 and 30203 by employing the Euclidean algorithm. Utilize prime factorization, considering the Fundamental Theorem of Arithmetic, and verify if both numbers are divisible by the same prime factors.

**Original Instruction 5:** Find the remainder when  $2 \times 12 \times 22 \times 32 \times \dots \times 72 \times 82 \times 92$  is divided by 5.

**Evol Instruction 5:** First, let  $P$  represent the product of the series, which can be expressed as  $P = \prod_{n=1}^9 (2 + 10n)$ . Next, calculate the value of  $P$ . Then, determine the remainder, denoted as  $R$ , when  $P$  is divided by 5. Ensure that  $R$  is a positive integer.

**Original Instruction 6:** Is the function  $f(x) = \lfloor x \rfloor + \frac{1}{2}$  even, odd, or neither? Enter odd, even, or neither.

**Evol Instruction 6:** Determine if the function  $f(x) = \lfloor x \rfloor + \frac{1}{2}$  exhibits parity (evenness or oddness) or neither, considering the mathematical definitions of even and odd functions. If  $x > 0$ , introduce a variable  $y$  and compare  $f(x)$  with  $g(y) = y^2$ . Provide a brief explanation for your answer. Enter f(x) is even, f(x) is odd, or f(x) is neither even nor odd.



1296 A.2 IRM PROMPT  
1297  
1298  
12991300 **Example 5: Instruction Quality Ranking Prompt**

1301 You are a senior mathematics grading teacher in university, very skilled in high difficulty fields such as  
1302 Intermediate Algebra, Precalculus, Prealgebra, Number Theory, Geometry, Counting & Probability, Algebra  
1303 and so on.

1304 Your task is to act as an impartial judge to evaluate the quality of math problems based on their definition  
1305 completeness and difficulty and rank a set of maths problems according to these criteria. Make sure that your  
1306 assessment takes into account the following rules:

1307  
1308 **1.\*\* Problem statement completeness and correctness:\*\***

- 1309 • Assess the clarity and accuracy of the definition of each math problem. Ensure that the problem  
1310 statement provides sufficient information, conditions, and constraints.
- 1311 • Consider whether the problem allows for multiple interpretations or if further clarification is needed.
- 1312 • Evaluate the clarity of mathematical notation and terminology used in the problem.

1313  
1314 **2.\*\*Conceptual difficulty:\*\***

- 1315 • Evaluates the complexity of each mathematical problem in terms of the underlying concepts  
1316 involved. Ensure a solid and sound understanding of the underlying principles, or advanced  
1317 mathematical concepts.
- 1318 • Consider the depth of mathematical knowledge required to address and solve each problem.
- 1319 • Assess whether the problem encourages critical thinking and the application of mathematical  
1320 principles.

1321 **3.\*\*Computational complexity:\*\***

- 1322 • Examine the computational complexity of each problem. Judge whether it involves complex  
1323 calculations, algebraic operations, or non-trivial numerical operations.
- 1324 • Consider whether the problem requires sophisticated computational techniques or algorithms or  
1325 whether it can be answered with existing mathematical knowledge.

1326 **4.\*\* Problem contextualisation:\*\***

- 1327 • Consider the relevance of each mathematical problem in the given context or practical application.  
1328 as well as being relevant or having a meaningful meaning in the practical context.
- 1329 • Evaluate whether the theory of the mathematical problem is detached from the facts, spurious, and  
1330 non-existent.

1331 Avoid any position biases and ensure that the order in which the math problems were presented does not  
1332 influence your decision.

1333 Do not allow the length of the problems to influence your evaluation.

1334 Do not favor certain mathematical theory of the problems. Be as objective as possible.

1335 Below is a list of a set of math problems that you need to rank according to the rules above from  
1336 most complete and clear (1) to least complete and clear (N) based on the comprehensiveness and difficulty  
1337 level of the maths problem. Also, consider the difficulty level from most challenging (1) to easiest (N). Your  
1338 output needs to be placed in the <Rank> </Rank>section.

1339 And Your output is in JSON list format, where each element is a dictionary with three keys:

- 1340 • instruction: represents the math problem.
- 1341 • score: represents the result of your ranking for the problem.
- 1342 • reason: provide your explanation in detail for your ranking result.

1343 [### Math Problems List ###]: <PROBLEMS\_HERE >  
1344  
1345  
1346  
1347  
1348  
1349

Math instructions	
1350	
1351	
1352	<b>Instruction 1.</b> Write $x^{10} + x^5 + 1$ as the product of two polynomials with integer coefficients. <b>Rank: 2.</b> Reason: The problem is well-defined and clear. It requires a deep understanding of polynomial factorization, making it conceptually difficult. The computational complexity is high as it involves the factorization of a high-degree polynomial.
1353	
1354	<b>Instruction 2.</b> Express the following sum as a common fraction:
1355	$\frac{1}{1 \cdot 2} + \frac{1}{2 \cdot 3} + \frac{1}{3 \cdot 4} + \frac{1}{4 \cdot 5} + \cdots + \frac{1}{9 \cdot 10}.$
1356	
1357	<b>Rank: 1.</b> Reason: This problem is clear and complete, with a well-defined mathematical notation. It requires a good understanding of series and fractions, making it conceptually challenging. The computational complexity is moderate as it involves the summation of fractions.
1358	
1359	
1360	
1361	GPT-4 Ranking
1362	<b>Instruction 3.</b> Given that I read 6 to 8 novels monthly, each novel having 200 to 350 pages, and considering a year may have 12 or 13 months due to a leap year or a year when I choose to read more, also taking into account that my reading speed can fluctuate between 1 to 2 novels per day, and the number of days in a month can be 28, 30, or 31, and the genre of the novel can influence the number of pages (for instance, a thriller may contain 300 to 350 pages, while a romance novel may contain 200 to 250 pages), calculate the total pages of novels I will read in a year. <b>Rank: 3.</b> Reason: This problem is complex and requires multiple steps to solve, making it conceptually and computationally challenging. It is well-contextualized and relevant, but the problem statement could be clearer and more concise.
1363	
1364	
1365	<b>Instruction 4.</b> James decides to run 3 sprints 3 times a week. He runs 60 meters each sprint. How many total meters does he run a week? <b>Rank: 4.</b> Reason: This problem is clear, complete, and well-contextualized. However, it is conceptually and computationally less challenging than the others, as it only requires basic multiplication.
1366	
1367	
1368	<b>Instruction 5.</b> What is $\sqrt{2}$ times $\sqrt{3}$ ? <b>Rank: 5.</b> Reason: This problem is clear and complete, but it lacks a practical context. It is conceptually and computationally less challenging than the others, as it only requires knowledge of the multiplication of square roots.
1369	
1370	
1371	<b>Instruction 6.</b> What is $\sqrt{-2}$ times $\sqrt{3}$ ? <b>Rank: 6.</b> Reason: This problem is clear and complete, but it lacks a practical context. It is conceptually challenging as it involves the multiplication of a real number and an imaginary number. However, the computational complexity is low.
1372	
1373	

Table 10: A case study GPT-4 scoring the evolved instructions from two aspects: i) Difficulty, and ii) Definition.

1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

1404 A.3 PRM PROMPT  
 1405

1406 **Example 6: PRM Step Level Labeling Prompt**  
 1407

1408 You are a senior mathematics grading teacher in university, very skilled in high difficulty fields such as  
 1409 Intermediate Algebra, Precalculus, Prealgebra, Number Theory, Geometry, Counting & Probability, Algebra  
 1410 and so on. Below is a mathematical problem and its corresponding solution, as well as a JSON list format for  
 the solution, where each element is a dictionary with two keys:

- 1411 • idx: represents the number of each step.
- 1412 • value: represents each step in the problem-solving process.

1413 Firstly please provide your judgement whether the solution is correct. Your judgment (which must be only  
 1414 True or False) needs to be placed in the <Judge> </Judge>section.

1415 And then you need to judge whether each step is correct and give a score for each solving step in the JSON  
 1416 list which needs to be placed in the <Scores> </Scores>section.

1417 There are three kinds of scores below:

- 1418 • 1: indicates that the step is correct.
- 1419 • 0: indicates that the step is ambiguity, meaningless, or subtly misleading, or not helpful to the  
 1420 entire problem-solving process.
- 1421 • -1: indicates that the step is incorrect.

1422 If this step leads to a final wrong answer, then rate -1. If not, rate 1 or 0.

1423 Here are some rules about whether the solution's each step is correct:

- 1424 • **## Problem-Solving Thoughts ##**: You should first think about how to solve this problem, and then  
 1425 judge whether this step is correct.
- 1426 • **## Calculation Accuracy ##**: You should carefully check and verify whether each step is calculated  
 1427 correctly, including various mathematical numerical calculations. Notably you don't need to  
 1428 consider simplification.
- 1429 • **## Logical Coherence ##**: You should judge whether each step is logically coherent and reasonable.
- 1430 • **## Basic Theories and Principles ##**: You should judge whether each step correctly is using basic  
 1431 mathematical theories, principles, or formulas.

1432 You need to constantly verify and check repeatedly whether each step is correct. And rate each step carefully,  
 1433 honestly, and without bias, order, or discrimination.

1434 Your output is a JSON list format, where each element is a dictionary with three keys:

- 1435 • idx: represents the number of each step.
- 1436 • score: represents your rating for this step, which can only be -1, 0 and 1.
- 1437 • reason: provide your explanation in detail for your rating whether each step of the problem-solving  
 1438 process is correct.

1439 [### Problem ###]: <INSTRUCTION\_HERE >  
 1440

1441  
 1442 [### Solution JSON List ###]: <Solution JSON List >  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

Question	<i>John and his two brothers decide to split the cost of an apartment. It is 40% more expensive than John's old apartment which costs \$1200 per month. How much does John save per year by splitting the apartment compared to living alone?</i>
	Step 1: Calculate the cost of John's old apartment. John's old apartment costs \$1200 per month. <b>Score: 1.</b>
	Step 2: Calculate the cost of the new apartment. <b>Score: 1.</b>
	Step 3: The new apartment costs 40% more than John's old apartment, so we need to find 40% of \$1200 and add that to the cost of the old apartment. <b>Score: 1.</b>
	Step 4: 40% of \$1200 is $(40/100) * \$1200 = \$480$ . So the new apartment costs $\$1200 + \$480 = \$1680$ per month. <b>Score: 1.</b>
	Step 5: Calculate the savings for John by splitting the apartment. <b>Score: 1.</b>
	Step 6: If John lived alone, he would pay \$1680 per month for the new apartment. <b>Score: 1.</b>
	Step 7: However, since he splits the cost with his two brothers, he only needs to pay 1/3 of the total cost. <b>Score: 1.</b>
GPT-4 Labeling	Step 8: So, John's share of the cost is $\$1680 / 3 = \$560$ per month. <b>Score: 1.</b>
	Step 9: Calculate the savings per year for John. <b>Score: 1.</b>
	<b>Step 10: There are 12 months in a year, so the total savings for John in a year would be <math>12 * \\$560 = \\$6720</math>. Score: -1.</b> Reason: The step is incorrect. The calculation is wrong. The total savings for John in a year should be calculated as the difference between what he would pay living alone and what he pays when splitting the cost, multiplied by 12. The correct calculation should be $(1200 - 560) * 12 = \$7680$ .
	<b>Step 11: So, John saves \$6720 per year by splitting the apartment compared to living alone. Score: -1.</b> Reason: The step is incorrect. The final answer is wrong. The correct answer should be \$7680.
	<b>Step 12: The answer is: 6720. Score: -1.</b> Reason: The step is incorrect. The final answer is wrong. The correct answer should be \$7680.

Table 11: A case study from Mistral-7B model on GSM8k training set. Red text denotes the incorrect steps that GPT-4 to successfully label errors.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

Question	Compute $\cos 120^\circ$ .
GPT-4 Labeling	Step 1: First, we need to find a right triangle with an angle of $120^\circ$ and then use the Pythagorean theorem to find the length of the hypotenuse. <b>Score: 0.</b>
	Step 2: The first step is to visualize a right triangle that has an angle of $120^\circ$ . <b>Score: -1.</b> Reason: This step is incorrect. A right triangle cannot have an angle of 120 degrees.
	Step 3: We can do this by drawing a line and then rotating it $120^\circ$ around the midpoint of the line. <b>Score: -1.</b> Reason: This step is incorrect. Rotating a line around its midpoint does not create a right triangle with an angle of 120 degrees.
	Step 4: This will give us a triangle with two $60^\circ$ angles and one $120^\circ$ angle. <b>Score: -1.</b> Reason: This step is incorrect. A triangle with two angles of 60 degrees and one angle of 120 degrees is not a right triangle.
	Step 5: Next, we need to find the length of the hypotenuse of this triangle. <b>Score: 0.</b>
	Step 6: We can do this by using the Pythagorean theorem, which states that in a right triangle, the square of the length of the hypotenuse is equal to the sum of the squares of the lengths of the other two sides. <b>Score: 0.</b>
	Step 7: In this case, the other two sides of the triangle are the two legs of the $60^\circ$ angles, which are each 1 unit long. <b>Score: -1.</b> Reason: This step is incorrect. The sides of a right triangle do not correspond to the legs of the 60 degree angles in this context.
	Step 8: So, we have $h^2 = 1 + 1 = 2$ . This means that the length of the hypotenuse $h$ is $\sqrt{2}$ . 'Now, we need to find the cosine of $120^\circ$ '. <b>Score: -1.</b> Reason: This step is incorrect. The calculation is based on incorrect assumptions from previous steps.
	Step 9: The cosine of an angle is defined as the ratio of the length of the side opposite the angle to the length of the hypotenuse. Now, we need to find the cosine of $120^\circ$ . <b>Score: -1.</b> Reason: This step is incorrect. The cosine of 120 degrees is not $1/\sqrt{2}$ .
	Step 10: The cosine of an angle is defined as the ratio of the length of the side opposite the angle to the length of the hypotenuse. <b>Score: 1.</b>
	Step 11: In this case, the opposite side is the leg of the $60^\circ$ angle, which is 1 unit long. So, we have $\cos 120^\circ = \frac{1}{\sqrt{2}}$ . <b>Score: -1.</b> Reason: This step is not entirely correct. The cosine of an angle is defined as the ratio of the adjacent side to the hypotenuse, not the opposite side to the hypotenuse.
	Step 12: Finally, we need to simplify this expression. The square root of 2 is usually written as $\sqrt{2}$ , so we can rewrite $\cos 120^\circ$ as $\frac{1}{\sqrt{2}}$ . <b>Score: -1.</b> Reason: This step is incorrect. The cosine of 120 degrees is not $1/\sqrt{2}$ .
	Step 13: So, the final answer is $\frac{1}{\sqrt{2}}$ . The answer is: $\frac{1}{\sqrt{2}}$ . <b>Score: -1.</b> Reason: This step is incorrect. The final answer is not $1/\sqrt{2}$ . The correct answer is $-1/2$ .

Table 12: A case study from Mistral-7B model on MATH training set. Red text denotes the incorrect steps that GPT-4 to successfully label errors.

#### A.4 THE EFFECTS OF MATH EVOL-INSTRUCT IN THE SFT AND RL STAGES.

##### 1. Math Evol-based Instruction works efficiently in the SFT stage.

We innovatively propose two evol instruct strategies for mathematical tasks: upward evolution and downward evolution. Through five rounds of iterative evolution, we successfully constructed a 210k dataset. In the supervised fine-tuning (SFT) phase, we conducted a comparative analysis between our dataset and manually annotated mathematical datasets, including GSM8k and MATH. In Table 13 results reveal that, with an equivalent dataset size of 15k, our evolved instructions yield improvements of 4.9% on GSM8k and 1.6% on MATH over human instructions. Furthermore, employing our total 210k evolved dataset leads to a 21.1% performance boost on GSM8k and a 15.2% enhancement on MATH. This substantiates the effectiveness of our mathematical instruction evolution strategies, significantly diminishing the dependence on laborious manual annotation efforts.

Table 13: Performance comparison on GSM8k and MATH using manually annotated GSM8k and MATH data and our math evol-instruct dataset. We employ the Mistral 7b model for training in the SFT stage.

Dataset	GSM8k	MATH
GSM8k and MATH, Human 15k	59.3	14.5
GSM8k and MATH, Human 7.5k + Evol 7.5k	62.9	15.2
GSM8k and MATH, Evol 15k	64.2	16.1
GSM8k and MATH, Human 15k + Evol 195k	80.4	29.7

##### 2. Math Evol-based Instruction works efficiently in the RL stage.

We propose the innovative use of math evolved instruction data in the reinforcement learning (RL) stage. In Table 14 we combine manually annotated mathematical data (i.e., GSM8k and MATH) with our evolved instructions data, utilizing IRM and PRM as reward models. The findings indicate that with a dataset size of 15k, our evolved instructions achieve a performance improvement of 0.5% on GSM8k and 0.6% on MATH in RL scenarios. Moreover, utilizing our comprehensive 210k evolved dataset results in performance gains of 3.3% on GSM8k and 3.9% on MATH during the RL stage. These outcomes significantly enhance RL model performance, indicating the effectiveness of our math-evolved instruction data in RL stage and addressing the scarcity of manually curated datasets.

Table 14: Performance comparison on GSM8k and MATH using manually annotated GSM8k and MATH data and our math evol-instruct dataset for training in the RL stage . We employ the WizardMath-Mistral 7b as our policy model.

Mistral-7B: WizardMath-SFT	GSM8k	MATH
GSM8k and MATH, Human 15k for RL	80.9	30.1
GSM8k and MATH, Human 7.5k + Evol 7.5k for RL	81.1	30.5
GSM8k and MATH, Evol 15k for RL	81.4	30.7
GSM8k and MATH, Human 15k + Evol 195k for RL	84.2	34.0

#### A.5 THE DIFFERENCE BETWEEN MATH EVOL-INSTRUCT AND WIZARDLM EVOL-INSTRUCT.

Inspired by prior studies, our math evol instruct method diverges from Wizardlm’s depth and breadth evolution strategies used for general tasks and WizardCoder for code task. We focus on both upward and downward evolution techniques for mathematical tasks, aiming to create a more complex and diverse math dataset. The Table 15 compares the performance between WizardLM’s original evol-instruct and our math-evol-instruct, the latter improves 7.8% on the GSM8k and 6.6% on the MATH. Therefore, it is imperative to utilize the math evol-instruct method specifically designed for mathematical applications rather than just leverage original evol-instruct. General evol-instruct strategies in WizardLM do not meet the requirements for the math scenario.

Although WizardLM and WizardCoder have also shown the effectiveness of evol-instruct in enhancing LLM’s instruction following ability, they just focus on the SFT stage. We not only verified and improved the effectiveness of math evol-instruct in SFT, but also firstly leveraged the evol-instruct in

reinforcement learning : the experimental results demonstrate the evolved data can further improve the performance of model in RL stage, thus we unlock the data limitations of reinforced math research.

Table 15: Performance comparison of WizardLM evol instruct and our WizardMath evol instruct on the GSM8k and MATH in SFT stage. The base model is Mistral 7b.

Dataset	GSM8k	MATH
WizardLM evol instruct, 210k	72.6	23.1
WizardMath evol instruct, 210k	80.4	29.7

#### A.6 DIFFERENCE BETWEEN CURRENT VERSION AND PREVIOUS VERSIONS.

The current version of WizardMath greatly improves mathematical reasoning skills. The enhancement in performance can be attributed to three main factors:

- Firstly, the data size was expanded from 96k to 210k entries.
- Secondly, the instruction evolution prompts were refined to be more precise and detailed, requiring GPT-4 to use a first plan then generate step-by-step approach for instruction evolution.
- Lastly, the hyper-parameters for RL training were further optimized, including adjustments to the learning rate, KL coefficients, and training steps.

The Table 16 below presents the performance improvements of the new model over the original one, utilizing 96k training data (a mix of evolved GSM8k and MATH data) on Llama2-7B: a 5.9% increase on GSM8k and a 5.7% on MATH during SFT. When including the RL, the overall gains rise to 6.8% on GSM8k and 6.1% on MATH.

Table 16: Explore the performance improvements of the new model over the original one in the SFT and RL.

Dataset	SFT		RL	
	GSM8k	MATH	GSM8k	MATH
WizardMath-original 96k	52.6	8.6	54.9	10.7
WizardMath-new 96k	58.5	14.3	61.7	16.8

#### A.7 EXPLORE GPT-4 PER STEP ACCURACY

We use GPT-4 to label randomly selected 200 samples from the PRM800k dataset and compared the results with Human-labeler, we employ F1 score as the metric to measure GPT-4 and manual annotations, we observed that GPT-4 exhibits 86% on the GSM8k and 72% consistency on the MATH with manual annotations, which indicates the effectiveness of our GPT-4 label method.

#### A.8 THE IMPACT OF THE DIFFERENT ROUND FOR UPWARD AND DOWNWARD EVOL-INSTRUCT

Table 18, Table 17 explore the impact of different rounds of upward and downward instruction evolution on GSM8k and MATH. By conducting 5 rounds of upward and downward evolution on GSM8k and MATH, and using the Mistral-7B base model for fine-tuning, we found that: Each round of upward evolution yielded a 7.5%-11.0% improvement on GSM8k and a 3.1%-7.2% increase on MATH over the baseline manual data. Similarly, each round of downward evolution demonstrated a 5.7%-10.3% improvement on GSM8k and a 2.1%-5.3% increase on MATH compared to the original manual data.

When merging the data from 3 rounds of upward evolution, we note peak performance for both GSM8k and MATH, with subsequent rounds leading to a gradual decline or plateau in performance. This phenomenon may be attributed to the instructions becoming more complex and abstract, which increases the rate of invalid instructions beyond GPT-4's capacity for accurate responses. Similarly,

when merging the data from 2 rounds of downward evolution, the performance also reaches the optimal level, but the performance slowly decreases after more than two rounds. This may be due to excessive evolution leading to instructions becoming too simple, lacking diversity, redundant, which increases the proportion of invalid data. Consequently, our study chosen three rounds of upward evolution and two rounds of downward evolution to balance diversity, complexity, and accuracy effectively.

Table 17: Impact of different Upward Evol-Instruct turns on Mistral-7B SFT.

Data	GSM8K							MATH						
	Original	1st	2nd	3rd	4th	5th	pass@1	Original	1st	2nd	3rd	4th	5th	pass@1
Original	✓	✗	✗	✗	✗	✗	59.7	✓	✗	✗	✗	✗	✗	15.1
Upward Evol	✓	✓	✗	✗	✗	✗	70.7	✓	✓	✗	✗	✗	✗	22.3
	✓	✗	✓	✗	✗	✗	70.6	✓	✗	✓	✗	✗	✗	21.8
	✓	✗	✗	✓	✗	✗	69.9	✓	✗	✗	✓	✗	✗	21.1
	✓	✗	✗	✗	✓	✗	68.1	✓	✗	✗	✗	✓	✗	19.7
	✓	✗	✗	✗	✗	✓	67.2	✓	✗	✗	✗	✗	✓	18.2
	✓	✓	✓	✗	✗	✗	73.4	✓	✓	✓	✗	✗	✗	24.1
	✓	✓	✓	✓	✗	✗	<b>75.9</b>	✓	✓	✓	✓	✗	✗	<b>25.6</b>
	✓	✓	✓	✓	✓	✗	75.5	✓	✓	✓	✓	✓	✗	25.3
	✓	✓	✓	✓	✓	✓	75.8	✓	✓	✓	✓	✓	✓	25.6

Table 18: Impact of different Downward Evol-Instruct turns on Mistral-7B SFT.

Data	GSM8K							MATH						
	Original	1st	2nd	3rd	4th	5th	pass@1	Original	1st	2nd	3rd	4th	5th	pass@1
Original	✓	✗	✗	✗	✗	✗	59.7	✓	✗	✗	✗	✗	✗	15.1
Downward Evol	✓	✓	✗	✗	✗	✗	70.0	✓	✓	✗	✗	✗	✗	20.4
	✓	✗	✓	✗	✗	✗	69.8	✓	✗	✓	✗	✗	✗	19.8
	✓	✗	✗	✓	✗	✗	68.5	✓	✗	✗	✓	✗	✗	19.1
	✓	✗	✗	✗	✓	✗	66.9	✓	✗	✗	✗	✓	✗	18.3
	✓	✗	✗	✗	✗	✓	65.4	✓	✗	✗	✗	✗	✓	17.2
	✓	✓	✓	✗	✗	✗	<b>72.1</b>	✓	✓	✓	✗	✗	✗	<b>23.1</b>
	✓	✓	✓	✓	✗	✗	72.1	✓	✓	✓	✓	✗	✗	22.8
	✓	✓	✓	✓	✓	✗	71.9	✓	✓	✓	✓	✓	✗	22.5
	✓	✓	✓	✓	✓	✓	71.8	✓	✓	✓	✓	✓	✓	23.0

#### A.9 EXPLORE THE IMPACT OF THE SCORING AGGREGATION STRATEGY AT EACH STEP OF THE PRM FOR THE RL POLICY TRAINING

Identifying incorrect steps is critical in the step-by-step math problem solving process. Even if the solution process is mostly correct, a single incorrect step often leads to an incorrect final answer and thus cannot be based on the number of correct steps. Our aim is to supervise the RL process by identifying the most error-prone steps, specifically those with the minimum reward scores.

In Table 19, we explore the impact of five score aggregation strategies on RL training:

- Max: uses the maximum of all step scores.
- Mean: calculates the mean of all step scores.
- Product: calculates the product of each step score.
- Dense per-step: considers each step scores to supervise RL training.
- Min: picks the minimum of all step scores.

The experimental results on Mistral-7B show that the Min strategy is the most effective, outperforming the Dense per-step strategy by 1.3% on GSM8k and 1.5% on MATH. The Max strategy is the worst, it's due to Max focus on the maximum reward scores for training, which leads to overlook incorrect supervisory signals and only the correct steps are mainly reinforced. And Dense per-step reward



strategy will diminish supervisory signals for erroneous steps, so we finally use the minimum reward at the end of the sequence.

Table 19: Explore the impact of the scoring aggregation strategy at each step of the PRM for the RL policy training. We use the WizardMath-Mistral-7B-SFT base policy model for the RL training.

Scoring Strategies	Max	Mean	Product	Dense per-step	Min
GSM8k	78.3	82.1	83.6	82.9	84.2
MATH	27.6	32.4	33.7	32.5	34.0

#### A.10 EXPLORE THE PERFORMANCE ON OUT OF DOMAIN (I.E,MATH) WHEN TRAINING ONLY GSM8K, AND VICE VERSA.

In Table 20 We only train SFT model on GSM8k data and evaluate out-of-domain performance (i.e., on MATH), and our WizardMath-GSM8k model attains a 8.1% accuracy, surpassing the manually annotated Human-GSM8k by 5.1%, respectively. Conversely, only training on MATH data and assessing out-of-domain performance (i.e., on GSM8k), WizardMath-MATH achieves a 39.4% accuracy, outperforming Human-MATH (13.8%) by 25.6%. These findings emphasize that our method significantly enhances performance on out-of-domain tasks.

Table 20: Explore the performance on out of domain (i.e,MATH) when training only GSM8k, and vice versa. We use the llama2 7b model for the SFT training.

Dataset	GSM8k	MATH
Human-GSM8k	41.6	3.0
WizardMath-GSM8k	61.9	8.1
Human-MATH	13.8	4.7
WizardMath-MATH	39.4	20.6
WizardMath	64.2	22.1

### 1782 A.11 DATA CONTAMINATION CHECK

1783  
1784 Apart from the performance analysis, we also investigate whether evolution leads to the data contami-  
1785 nation between training data and test set. To address this consideration, we employ instructions in the  
1786 GSM8k and MATH test set as queries to retrieve the top-5 samples from all evolved training data  
1787 with an embedding model, gte-large (Li et al., 2023d). Additionally, we employ GPT-4 to provide  
1788 similarity judgement between the test sets and the retrieved samples, and remove the top-2 similar in-  
1789 structions. The prompt and details are shown in Appendix A.12. Figure 4 in Appendix illustrates that  
1790 the evolution process does not yield higher similarity scores. Furthermore, similarity scores across all  
1791 rounds remain relatively low. These findings indicate that the primary source of performance gain  
1792 is the introduction of more complex and comprehensive data based on our downward and upward  
1793 instruction evolution.

### 1794 A.12 SIMILARITY CHECKING AND DATA FILTERING

1795  
1796 The prompt formats to compute the similarity score between two given math problem tasks are as  
1797 follow:

#### 1798 Example 7: System Prompt for Similarity Checking

1800 Your task is to evaluate the similarity of the two given math problems. Please review the two math problem  
1801 tasks carefully, paying close attention to the overlap in variables, conditions, participants, actions, or processes,  
1802 topics, and contents and core concept and structure. Once you have carefully reviewed both math problem  
1803 tasks, provide a similarity score between these two math problem tasks. The score should range from 1 to 10  
1804 (1: completely different math problem tasks; 10: identical math problem tasks). You only need to provide  
1805 your score without any explanation.

#### 1806 # Problem-1

1807 {task1}

#### 1808 # Problem-2

1809 {task2}

1810 Your judgement score:

1811  
1812 To thoroughly prevent data leakage from the GSM8k and MATH test datasets to the training dataset,  
1813 we implemented an additional data filtering step. Utilizing the SOTA embeddings model, gte-large,  
1814 we treated all test samples as queries to extract the top 5 samples from the training data. Following  
1815 this, GPT-4 was employed to evaluate the similarity between the retrieved samples and the test set.

### 1816 A.13 DETAIL WORKS

## 1817 B RELATED WORK

1820 **Large Language Models.** LLMs have significantly advanced Natural Language Processing, with  
1821 models like OpenAI’s GPT Series (Brown et al., 2020a; OpenAI, 2023), Anthropic’s Claude (Bai  
1822 et al., 2022), Google’s PaLM (Chowdhery et al., 2022; Anil et al., 2023), Gemini (Team, 2023), and  
1823 Gemma (Team et al., 2024) featuring billions of parameters and trained on massive textual datasets.  
1824 The AI field has also seen a rise in open-source LLMs such as Mistral (Jiang et al., 2023), Llama  
1825 Series (Touvron et al., 2023a;b; Dubey et al., 2024; Taylor et al., 2022), DeepSeek (Bi et al., 2024;  
1826 Shao et al., 2024), Qwen (Bai et al., 2023; Yang et al., 2024) etc. (Zeng et al., 2022; Penedo et al.,  
1827 2023; Scao et al., 2022). Notably, Llama serves as a foundational model for supervised fine-tuning,  
1828 leading to the development of models like Alpaca, Vicuna, Guanaco, and Orca (Taori et al., 2023;  
1829 Chiang et al., 2023; Dettmers et al., 2023; Mukherjee et al., 2023).

1830 **Large Language Models For Mathematical reasoning.** NLP models face challenges with complex  
1831 reasoning, including mathematical (Lu et al., 2022; Frieder et al., 2023; Long et al., 2024; Zhang et al.,  
1832 2024c; Xia et al., 2024), common-sense (Talmor et al., 2019; Geva et al., 2021). Significant research  
1833 focuses on Mathematical Word Problems (MWP), which demand understanding of mathematical  
1834 concepts and multi-step reasoning (Koncel-Kedziorski et al., 2016; Patel et al., 2021; Lan et al.,  
1835 2022; Cobbe et al., 2021; Jie et al., 2022; Yuan et al., 2023b; Fu et al., 2023b; Zheng et al., 2023;  
Zhao et al., 2023; Wang et al., 2023b; Imani et al., 2023; Yuan et al., 2023a; Wang et al., 2023e;

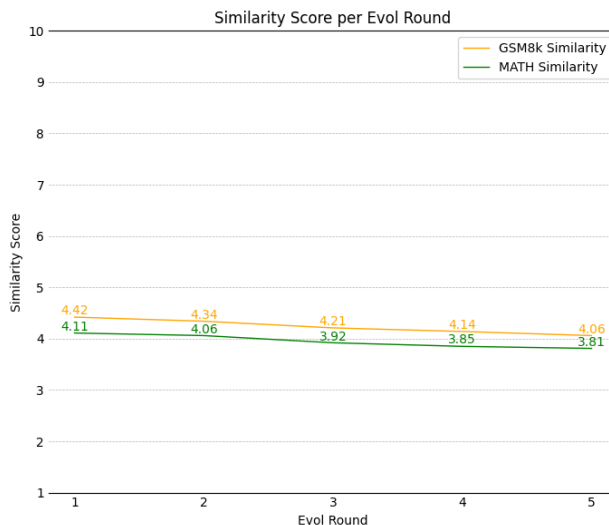


Figure 4: Average similarity scores between GSM8k, MATH samples and the top-1 retrieved data for each round.

Ahn et al., 2024). Models are tested on various MWP benchmarks (Hosseini et al., 2014; Roy & Roth, 2015; Koncel-Kedziorski et al., 2015; Patel et al., 2021; Cobbe et al., 2021; Ling et al., 2017; Hendrycks et al., 2021; Xia et al., 2024; Huang et al., 2024a; Zhang et al., 2024c; Li et al., 2024d; Anand et al., 2024). Techniques like Chain-of-Thought Prompting (Wei et al., 2022), Least-to-Most prompting (Zhou et al., 2022), and Complex CoT (Fu et al., 2022) enhance reasoning by introducing multiple steps and breaking problems into sub-problems. There are some models aimed at improving math CoT reasoning skills such as MetaMath (Yu et al., 2023b), MathScale (Tang et al., 2024), Xwin-math (Li et al., 2024a) etc. (Liu et al., 2023; Liu & Yao, 2024; Jiang et al., 2024b; Xiong et al., 2024; Chan et al., 2024; Lin et al., 2024; Huang et al., 2024b; Mitra et al., 2024b; Yuan et al., 2023a; Wu et al., 2024; Mitra et al., 2024a; Zhang et al., 2024e; Fu et al., 2023a; Ni et al., 2024; Tian et al., 2024; Zhang et al., 2024b; Zeng et al., 2024; Shao et al., 2024; Feng et al., 2023; Wu et al., 2023a; Tong et al., 2024; Yang et al., 2024; Zhou et al., 2024; Chen et al., 2024c; Zhang et al., 2024d; Su et al., 2024). Some models enhance mathematical reasoning by integrating python tools, such as TORA (Gou et al., 2023), MAMmoTH (Yue et al., 2023), Openmathinstruct (Toshniwal et al., 2024), etc. (Lu et al., 2024; Li et al., 2024c; Yu et al., 2024; Wang et al., 2023a; Li et al., 2024b)

**Large Language Models For Reinforcement Learning.** State-of-the-art models often display logical errors and illusions, particularly in domains requiring complex, multi-step reasoning, leading to significant challenges (Bubeck et al., 2023; Maynez et al., 2020). Strategies such as training reward models help discriminate between desirable and undesirable outputs (Lightman et al., 2023; Wu et al., 2023b; Chen et al., 2024b). Historically, outcome-based approaches focused on algorithmic tasks (Graves et al., 2014; Reed & De Freitas, 2015; Li et al., 2016; Cai et al., 2017; Yu et al., 2023a), while recent research demonstrates the efficacy of reward models or validators in enhancing model performance (Cobbe et al., 2021; Wang et al., 2023c;d; Li et al., 2022a;b). Reward models have also been incorporated into reinforcement learning pipelines and employed in rejection sampling to align Large Language Models (LLMs) with human preferences (Ziegler et al., 2019; Stiennon et al., 2020; Nakano et al., 2021; Ouyang et al., 2022; Nichols et al., 2020; Shen et al., 2021; Bai et al., 2022; Yuan et al., 2023c; Dong et al., 2023; Song et al., 2023; Touvron et al., 2023b). A contrast is drawn between outcome-supervised and process-supervised reward models, with the latter being more effective at addressing discrepancies arising from incorrect reasoning paths leading to correct outcomes (Uesato et al., 2022; Zelikman et al., 2022; Creswell et al., 2022). Recent advances have promoted process-based supervision through manual annotation, significantly benefiting LLMs over outcome-based approaches (Lightman et al., 2023; Zhu et al., 2023; Ni et al., 2022; Wang et al., 2024a; Sun et al., 2024; Chen et al., 2024a; Wang et al., 2024b; Zhang et al., 2024a;b). In this paper, we leverage AI models like ChatGPT to automatically offer annotation to improve the efficiency of this research line.

## C WIZARDMATH REBUTTAL

### C.1 REVIEWER-8CQG

Dear Reviewer 8CQg,

we thank you for your valuable comments and the time you spent reviewing our work! Your professional feedback provides valuable guidance for writing a more comprehensive and competitive paper. Below, we provide detailed responses to the **Weaknesses** and **Questions** raised in your review of our paper, addressing each point systematically.

Meanwhile, **in Appendix C.1 of our latest upload of revised paper (pages 36–49, lines 1892–2589)**, we also have added the discussions with the Reviewer-8CQg on the weaknesses and questions of our paper to respond to the Reviewer-8CQg’s comments and to further improve the research work.

#### C.1.1 WEAKNESSES-1

**The primary concern with this paper is the unfair comparison of baseline models in the results. While the authors claim that both supervised fine-tuning (SFT) with Math Evol-Instruct and reinforcement learning (RL) with the Instruction Reward Model (IRM) and Process Reward Model (PRM) are beneficial for enhancing mathematical reasoning, these approaches—SFT with synthesized data and the use of various reward models for RL—represent parallel research lines.**

We sincerely appreciate your attention to our work and your careful and responsible review and thank you for your valuable suggestions. To ensure a fair comparison, we conducted evaluations using WizardMath-SFT against all current state-of-the-art (SOTA) models across different scales of base models, as presented **in Table 21 and Table 22, Appendix C.1.2 of our latest upload of revised paper (pages 38–39, lines 2009–2105)**. The results confirm the effectiveness of our proposed Math Evol-Instruct approach. Meanwhile, during the PPO training stage, we applied IRM and PRM to different SFT backbones, significantly enhancing the mathematical reasoning ability of these models. This demonstrates the effectiveness and generalizability of our IRM and PRM methods refer to Weaknesses1.1-Weaknesses1.4 below for details.

Below, we provide detailed responses to Weaknesses 1.1- Weaknesses1.4 in the order you were raised.

#### C.1.2 WEAKNESSES-1.1

**In Table 1, the authors compare their model, which has undergone both SFT and RL, with models that have only undergone SFT. This comparison is unfair because these SFT models could also be further enhanced with RL techniques to improve mathematical reasoning (e.g., using ORM for RL on DartMath). It would be more appropriate to isolate the effects of SFT and RL for a fair comparison. In Table 1, the authors should compare the performance of models that have undergone SFT with Math Evol-Instruct against existing baselines such as MetaMath and DartMath. Additionally, comparisons with baselines like MetaMath and DartMath on the LLaMA-3.2 backbone would be valuable, as their training data is publicly available.**

We sincerely appreciate your insightful questions and detailed observations. To provide a more comprehensive and fair comparison, we have included the WizardMath-SFT results **in Table 21 and Table 22, Appendix C.1.2 of our latest upload of revised paper (pages 38–39, lines 2009–2105)**. These results evaluate the performance of **WizardMath-SFT**, trained exclusively using SFT, against current SOTA models across various base models. The key findings are summarized as follows:

#### 1. Performance Comparison:

- On **Llama-2-7B** and **Mistral-7B-v0.1**, WizardMath-SFT performs marginally below SOTA models (i.e., Xwin-Math and Skywork-Math) and outperforms existing other excellent models (i.e., DART-Math).
- On **Llama-2-13B** and **Llama-2-70B**, WizardMath-SFT achieves comparable performance to Xwin-Math.

- 1944
- On all various base models, WizardMath-SFT surpasses most existing SOTA models trained solely with SFT(i.e., DART-Math).

1945

1946

1947

1948

1949

Notably, WizardMath-SFT achieves these results using only 418K synthetic data points, a significantly smaller dataset compared to DART-Math (580k-590k), Xwin-Math (1440K) and Skywork-Math (2500K).

## 2. Comparison with advanced data synthesis methods (i.e., DART-Math, MetaMath)

1950

1951

1952

1953

1954

As shown in the following Table 23, DART-Math demonstrates strong performance across various base models and the data synthesis method proposed by DART-Math shows the effectiveness and outstanding performance. Meanwhile, WizardMath-SFT demonstrates comparable or superior performance to advanced data synthesis methods, such as **DART-Math** and **MetaMath**, across all base models. Key observations include:

- On **Mistral-7B-v0.1** and **DeepSeekMath**, WizardMath-SFT performs on par with DART-Math (Uniform & Prop2Diff) on GSM8k and surpasses DART-Math (Uniform & Prop2Diff) on MATH;
- On **Llama3.2 1B**, **Llama3.2 3B**, **Llama3-8B**, and **Llama3.1-8B**, **Llama2-7B**, WizardMath-SFT exhibits a 2%–7% improvement over DART-Math (Uniform & Prop2Diff) on the GSM8k benchmark. On the **MATH** benchmark, WizardMath-SFT outperforms DART-Math (Uniform & Prop2Diff) by approximately 5% – 10%.

1955

1956

1957

1958

1959

1960

1961

1962

1963

1964

These findings highlight the effectiveness of the proposed **Math Evol-Instruct** for enhancing mathematical reasoning capabilities.

1965

1966

1967

1968

1969

Notably, to ensure the same training settings as in our paper during the SFT stage, we employ a learning rate of  $2e-5$  for the Llama series base models (i.e., Llama2 7B, Llama3.1 8B, Llama3.2 1B, and Llama3.2 3B) and a learning rate of  $5e-6$  for Mistral-7B-v0.1. All models are trained for 3 epochs with a batch size of 256, and 4 checkpoints are saved per epoch. Finally, we select the checkpoint with the highest accuracy on the GSM8k and MATH benchmarks for reporting.

1970

1971

1972

1973

1974

1975

1976

1977

1978

1979

1980

1981

1982

1983

1984

1985

1986

1987

1988

1989

1990

1991

1992

1993

1994

1995

1996

1997

We have added the discussions about the Weaknesses-1.1 in **Appendix C.1.2 of our latest upload of revised paper (pages 36–40, lines 1921–2135)**

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

Table 21: In the study, we compare the WizardMath-SFT/RL model across various base models (0.1B-3B) with the SOTA models on the GSM8k and Math benchmarks. We report the Chain of Thought (CoT) pass@1 results without using any external Python tools. The results from 7B to 70B are shown in Table 22.

Model	Base	Params	GSM8k	MATH
Proprietary models				
GPT-o1 (OpenAI, 2023)	-	-	-	94.8
GPT-o1-mini (OpenAI, 2023)	-	-	-	90.0
Gemini-1.5 002 (Team et al., 2023)	-	-	-	86.5
Claude 3.5 Sonnet (Bai et al., 2022)	-	-	96.4	71.1
GPT-4o-2024-0513 (OpenAI, 2023)	-	-	96.1	76.6
GPT-4-turbo-0125 (OpenAI, 2023)	-	-	94.2	64.5
GPT-4-0314 (OpenAI, 2023)	-	-	94.7	52.6
GPT-4 (original version) (OpenAI, 2023)	-	-	92.0	42.5
Baichuan-3 (Yang et al., 2023)	-	-	88.2	49.2
GLM-4 (GLM et al., 2024)	-	-	87.6	47.9
Gemini Pro (Team, 2023)	-	-	86.5	32.6
Claude2 (Bai et al., 2022)	-	-	85.2	32.5
GPT-3.5-Turbo (OpenAI, 2023)	-	-	81.6	43.1
PaLM2 (Anil et al., 2023)	-	-	80.7	34.3
Minerva (Lewkowycz et al., 2022)	-	540B	58.8	33.6
GPT3.5 (Brown et al., 2020a)	-	-	57.1	-
Open-Source Models (0.1B-3B)				
GPT-2-Small (Brown et al., 2020b)	-	0.1B	6.9	5.4
GPT-2-Medium (Brown et al., 2020b)	-	0.3B	11.2	6.2
GPT-2-Large (Brown et al., 2020b)	-	0.7B	13.6	6.4
GPT-2-XL (Brown et al., 2020b)	-	1.5B	15.4	6.9
WizardMath-GPT-SFT	GPT-2-Small	0.1B	21.2	9.1
WizardMath-GPT-RL	GPT-2-Small	0.1B	26.4	12.3
WizardMath-GPT-SFT	GPT-2-Medium	0.3B	30.6	11.4
WizardMath-GPT-RL	GPT-2-Medium	0.3B	38.7	15.6
WizardMath-GPT-SFT	GPT-2-Large	0.7B	43.7	16.4
WizardMath-GPT-RL	GPT-2-Large	0.7B	50.1	21.2
WizardMath-GPT-SFT	GPT-2-XL	1.5B	51.9	18.3
WizardMath-GPT-RL	GPT-2-XL	1.5B	58.9	25.4
WizardMath-Qwen-SFT	Qwen-Math-2.5	1.5B	82.3	62.1
WizardMath-Qwen-RL	Qwen-Math-2.5	1.5B	86.7	68.6
Llama-3.2-Instruct (Dubey et al., 2024)	Llama 3.2	1B	44.4	30.6
MetaMath (Yu et al., 2023b)	Llama 3.2	1B	51.9	15.5
DART-Math-Prop2Diff (Tong et al., 2024)	Llama 3.2	1B	49.2	23.4
DART-Math-Uniform (Tong et al., 2024)	Llama 3.2	1B	55.8	22.0
WizardMath-Llama-SFT	Llama 3.2	1B	57.1	29.7
WizardMath-Llama-RL	Llama 3.2	1B	63.3	33.5
Llama-3.2-Instruct (Dubey et al., 2024)	Llama 3.2	3B	77.7	48.0
MetaMath (Yu et al., 2023b)	Llama 3.2	3B	72.6	25.9
DART-Math-Prop2Diff (Tong et al., 2024)	Llama 3.2	3B	74.0	37.8
DART-Math-Uniform (Tong et al., 2024)	Llama 3.2	3B	77.8	36.4
WizardMath-Llama-SFT	Llama 3.2	3B	80.3	45.2
WizardMath-Llama-RL	Llama 3.2	3B	85.5	49.9

Table 22: Continue Table 21, in this study, we compare the WizardMath-SFT/RL model across various base models (7B-70B) with the SOTA models on the GSM8k and Math benchmarks. We report the Chain of Thought (CoT) pass@1 results without using any external Python tools.

Model	Base	Params	GSM8k	MATH
Open-Source Models (7B-8B)				
Llama-2 (Touvron et al., 2023b)	-	7B	14.6	2.5
MAMmoTH-CoT (Yue et al., 2023)	Llama-2	7B	50.5	10.4
MathScale (Tang et al., 2024)	Llama-2	7B	66.3	31.1
MetaMath (Yu et al., 2023b)	Llama-2	7B	66.5	19.8
MuggleMath (Li et al., 2023a)	Llama-2	7B	68.4	-
Skywork-Math (Zeng et al., 2024)	Llama-2	7B	72.9	47.7
Math-Shepherd (Wang et al., 2024a)	Llama-2	7B	73.2	21.6
DART-Math-Prop2Diff (Tong et al., 2024)	Llama-2	7B	69.9	30.7
DART-Math-Uniform (Tong et al., 2024)	Llama-2	7B	73.8	29.5
Xwin-Math (Li et al., 2024a)	Llama-2	7B	82.6	40.6
WizardMath-Llama-SFT	Llama-2	7B	77.4	35.6
WizardMath-Llama-RL	Llama-2	7B	84.1	43.5
-----				
Mistral-v0.1 (Jiang et al., 2023)	-	7B	42.9	12.9
MathScale (Tang et al., 2024)	Mistral-v0.1	7B	74.8	35.2
MMIQc (Liu & Yao, 2024)	Mistral-v0.1	7B	74.8	36.0
MetaMath (Yu et al., 2023b)	Mistral-v0.1	7B	77.9	28.6
DART-Math-Prop2Diff (Tong et al., 2024)	Mistral-v0.1	7B	81.1	45.5
KPMath-Plus (Huang et al., 2024b)	Mistral-v0.1	7B	82.1	46.8
DART-Math-Uniform (Tong et al., 2024)	Mistral-v0.1	7B	82.6	43.5
Skywork-Math (Zeng et al., 2024)	Mistral-v0.1	7B	83.9	51.2
Math-Shepherd (Wang et al., 2024a)	Mistral-v0.1	7B	84.1	33.0
MAMmoTH2-Plus (Yue et al., 2024)	Mistral-v0.1	7B	84.7	45.0
JiuZhang3.0 (Zhou et al., 2024)	Mistral-v0.1	7B	88.6	52.8
Xwin-Math (Li et al., 2024a)	Mistral-v0.1	7B	89.2	43.7
WizardMath-Mistral-SFT	Mistral-v0.1	7B	82.8	48.1
WizardMath-Mistral-RL	Mistral-v0.1	7B	90.7	55.4
WizardMath-Mistral-SFT	Mistral-v0.3	7B	84.5	49.9
WizardMath-Mistral-RL	Mistral-v0.3	7B	90.4	55.6
WizardMath-Mathstral-SFT	Mathstral-v0.1	7B	88.3	64.2
WizardMath-Mathstral-RL	Mathstral-v0.1	7B	93.8	70.9
-----				
Qwen2.5-Math-Base (Yang et al., 2024)	Qwen2.5-Math	7B	91.6	55.4
WizardMath-Qwen-SFT	Qwen2.5-Math	7B	92.3	72.3
WizardMath-Qwen-RL	Qwen2.5-Math	7B	93.9	77.8
WizardMath-Qwen-SFT	Qwen2.5	7B	89.8	68.1
WizardMath-Qwen-RL	Qwen2.5	7B	94.0	74.5
-----				
DeepSeekMath-Base (Shao et al., 2024)	-	7B	64.2	36.2
NuminaMath-CoT (Li et al., 2024c)	DeepSeekMath	7B	75.4	55.2
MMIQc (Liu & Yao, 2024)	DeepSeekMath	7B	79.0	45.3
KPMath-Plus (Huang et al., 2024b)	DeepSeekMath	7B	83.9	48.8
DART-Math-Prop2Diff (Tong et al., 2024)	DeepSeekMath	7B	86.8	53.6
DeepSeekMath-RL (Shao et al., 2024)	DeepSeekMath	7B	88.2	51.7
DART-Math-Uniform (Tong et al., 2024)	DeepSeekMath	7B	88.2	52.9
WizardMath-DeepSeek-SFT	DeepSeekMath	7B	88.9	58.2
WizardMath-DeepSeek-RL	DeepSeekMath	7B	91.0	64.6
-----				
MetaMath (Yu et al., 2023b)	Llama 3	8B	77.3	20.6
MMIQc (Liu & Yao, 2024)	Llama 3	8B	77.6	29.5
DART-Math-Prop2Diff (Tong et al., 2024)	Llama 3	8B	81.1	46.6
DART-Math-Uniform (Tong et al., 2024)	Llama 3	8B	82.5	45.3
MAMmoTH2-Plus (Yue et al., 2024)	Llama 3	8B	84.1	42.8
Llama 3.1-Instruct (Dubey et al., 2024)	Llama 3	8B	84.5	51.9
JiuZhang3.0 (Zhou et al., 2024)	Llama 3	8B	88.6	51.0
WizardMath-Llama-SFT	Llama 3	8B	88.9	53.3
WizardMath-Llama-RL	Llama 3	8B	90.3	58.8
-----				
MetaMath (Yu et al., 2023b)	Llama 3.1	8B	80.4	35.4
DART-Math-Prop2Diff (Tong et al., 2024)	Llama 3.1	8B	84.3	46.5
DART-Math-Uniform (Tong et al., 2024)	Llama 3.1	8B	86.7	45.1
WizardMath-Llama-SFT	Llama 3.1	8B	89.2	55.8
WizardMath-Llama-RL	Llama 3.1	8B	93.4	62.3
-----				
Open-Source Models (13B)				
Llama-2 (Touvron et al., 2023b)	-	13B	28.7	3.9
MAMmoTH-CoT (Yue et al., 2023)	Llama 2	13B	56.3	12.9
MathScale (Tang et al., 2024)	Llama 2	13B	71.3	33.8
MetaMath (Yu et al., 2023b)	Llama 2	13B	72.3	22.4
MuggleMath (Li et al., 2023a)	Llama 2	13B	74.0	-
KPMath-Plus (Huang et al., 2024b)	Llama 2	13B	81.6	41.0
Xwin-Math (Li et al., 2024a)	Llama 2	13B	88.1	44.9
WizardMath-Llama-SFT	Llama 2	13B	86.8	46.5
WizardMath-Llama-RL	Llama 2	13B	89.7	50.6
-----				
Open-Source Models (70B)				
Llama-2 (Touvron et al., 2023b)	-	70B	56.8	13.5
MAMmoTH-CoT (Yue et al., 2023)	Llama-2	70B	72.4	21.1
MetaMath (Yu et al., 2023b)	Llama-2	70B	82.3	26.6
KPMath-Plus (Huang et al., 2024b)	Llama-2	70B	87.4	48.6
Xwin-Math (Li et al., 2024a)	Llama-2	70B	90.6	52.8
WizardMath-Llama-SFT	Llama-2	70B	89.5	54.4
WizardMath-Llama-RL	Llama-2	70B	92.8	58.6

Table 23: In this study, we mainly compare the performance of WizardMath-SFT with advanced data synthesis methods such as DART-Math and MetaMath on different base models under the GSM8k and MATH benchmarks in the SFT stage. We report the CoT pass@1 results of the model without relying on any external Python tools.

Model	Base	Params	GSM8k	MATH
DART-Math-Prop2Diff	Llama 3.2	1B	49.2	23.4
MetaMath	Llama 3.2	1B	51.9	15.5
DART-Math-Uniform	Llama 3.2	1B	55.8	22.0
WizardMath-Llama-SFT	Llama 3.2	1B	57.1	29.7
MetaMath	Llama 3.2	3B	72.6	25.9
DART-Math-Prop2Diff	Llama 3.2	3B	74.0	37.8
DART-Math-Uniform	Llama 3.2	3B	77.8	36.4
WizardMath-Llama-SFT	Llama 3.2	3B	80.3	45.2
MetaMath	Llama-2	7B	66.5	19.8
DART-Math-Prop2Diff	Llama-2	7B	69.9	30.7
DART-Math-Uniform	Llama-2	7B	73.8	29.5
WizardMath-Llama-SFT	Llama-2	7B	77.4	35.6
MetaMath	Mistral-v0.1	7B	77.9	28.6
DART-Math-Prop2Diff	Mistral-v0.1	7B	81.1	45.5
DART-Math-Uniform	Mistral-v0.1	7B	82.6	43.5
WizardMath-Mistral-SFT	Mistral-v0.1	7B	82.8	48.1
DART-Math-Prop2Diff	DeepSeekMath	7B	86.8	53.6
DART-Math-Uniform	DeepSeekMath	7B	88.2	52.9
WizardMath-DeepSeek-SFT	DeepSeekMath	7B	88.9	58.2
MetaMath	Llama 3	8B	77.3	20.6
DART-Math-Prop2Diff	Llama 3	8B	81.1	46.6
DART-Math-Uniform	Llama 3	8B	82.5	45.3
WizardMath-Llama-SFT	Llama 3	8B	88.9	53.3
MetaMath	Llama 3.1	8B	80.4	35.4
DART-Math-Prop2Diff	Llama 3.1	8B	84.3	46.5
DART-Math-Uniform	Llama 3.1	8B	86.7	45.1
WizardMath-Llama-SFT	Llama 3.1	8B	89.2	55.8

### C.1.3 WEAKNESSES-1.2

In analyzing the impact of training data size, the authors should compare their approach with the best method for SFT using synthesized data, specifically DartMath. MetaMath, which was developed around a year ago, uses GPT-3.5-turbo for data augmentation, making it an outdated and potentially unfair baseline.

Thank you for your insightful advice. In Appendix C.1.3 of our latest upload of revised paper (pages 40–41, lines 2137–2154), Figure 5, we explore the performance of WizardMath Evol-instruct in comparison with DART-Math and MetaMath across different training data scales on the GSM8k and MATH benchmarks in the SFT stage.

As the volume of training data increases, WizardMath-Evol-Instruct consistently improves its performance on the GSM8k and MATH benchmarks, exhibiting a slightly higher growth rate than DART-Math. In the initial stages, WizardMath slightly underperforms compared to DART-Math. This advantage may stem from DART-Math being distilled from DeepSeekMath-RL, an advanced mathematical reasoning model pre-trained on 120B high-quality mathematical tokens, showcasing exceptional proficiency in mathematical reasoning. However, once the dataset exceeds 60k, its performance begins to surpass DART-Math. At a data scale of 390k, WizardMath slightly outperforms DART-Math by 2%–3% on GSM8k and by 5%–6% on MATH. Additionally, WizardMath-Evol-Instruct consistently exceeds MetaMath at the same data scales, achieving increases of 3%–6% on GSM8k and 15%–20% on MATH. This performance gain is attributed to the efficiency of Math Evol-Instruct’s upward and downward evolution processes. These findings demonstrate that our Math Evol-Instruct method is also as scalable and effective as DART-Math for the large-scale synthetic data.



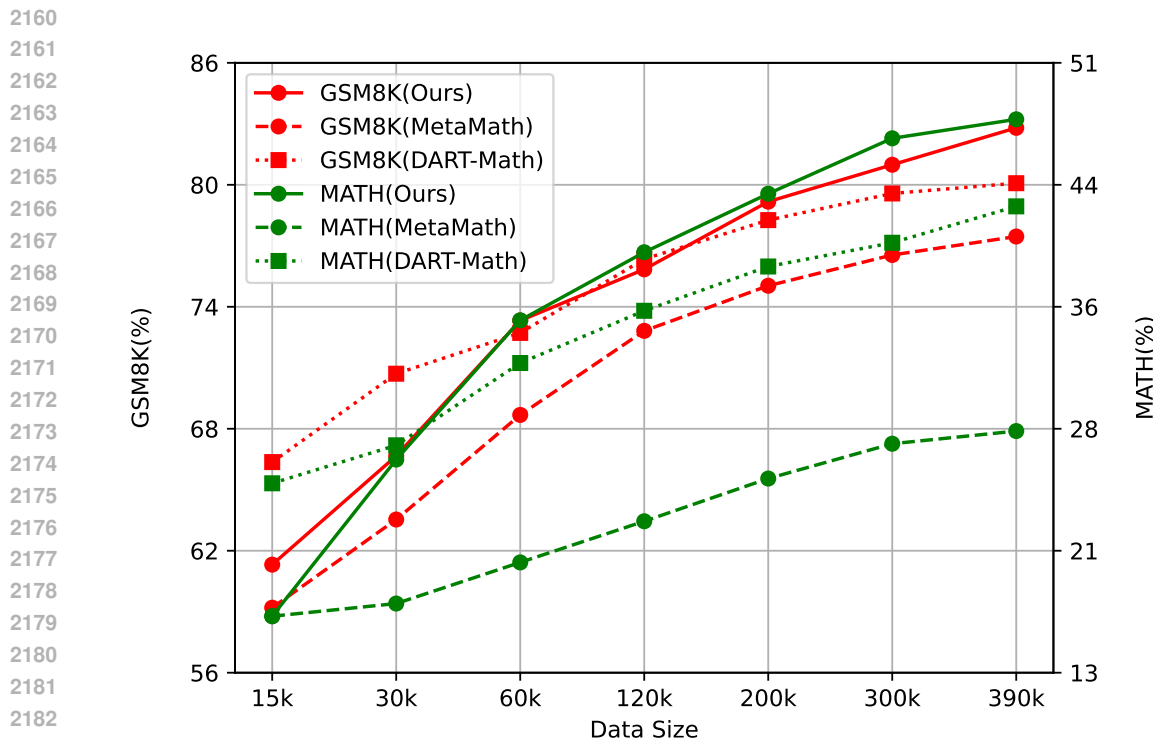


Figure 5: The performance of WizardMath Evol-instruct in comparison with DART-Math and MetaMath across different training data scales on the GSM8k and MATH benchmarks in the SFT stage. We use the Mistral-7B as base model

We have added the discussions about the **Weaknesses-1.2** in **Appendix C.1.3** of our latest upload of revised paper (pages 40-41, lines 2139–2189)

#### C.1.4 WEAKNESSES-1.3

It appears that SFT with Math Evol-Instruct yields inferior results compared to other SFT methods. From Table 4, the LLaMA2-7B: WizardMath-SFT scores 35.6 on MATH, which lags behind models like XwinMath and Skywork. Likely, it would also lag behind LLaMA2-7B fine-tuned on the DartMath training data. This suggests that the main contribution of the paper is in the RL component. Therefore, the primary focus should be on the results obtained with different reward models, as presented in Table 4, utilizing various SFT backbones.

Thank you for your valuable questions and insightful suggestions. Below are detailed responses to each question.

**Q1: It appears that SFT with Math Evol-Instruct yields inferior results compared to other SFT methods. From Table 4, the LLaMA2-7B: WizardMath-SFT scores 35.6 on MATH, which lags behind models like XwinMath and Skywork. Likely, it would also lag behind LLaMA2-7B fine-tuned on the DartMath training data.**

In Table 24 below, we show the performance comparison of WizardMath-SFT with DART-Math, Xwin-Math and Skywork-Math on the Llama2-7B base model on the MATH benchmark.

- **WizardMath-SFT vs. DART-Math:**

WizardMath-SFT, based on the Llama2-7B model, outperforms DART-Math-Uniform by 6.1% and DART-Math-Prop2Diff by 4.9% on the MATH. Notably, the amount of data used by WizardMath-SFT is only 70%–71% of DART-Math (418k vs. 591k; 418k vs. 585k).

- **WizardMath vs. Xwin-Math:**

Although WizardMath-SFT is 5% lower than Xwin-Math on the MATH, the amount of data used is only 29.0% of Xwin-Math (418k vs. 1440k), which is much less than Xwin-Math.

Moreover, Xwin-Math leverages GPT-4-turbo for data synthesis. However, WizardMath-SFT outperforms Xwin-Math on the MATH when using different backbones such as Mistral-7B-v0.1, Llama2-13B, and Llama2-70B as shown in **Table 22 of our latest upload of revised paper (pages 39, lines 2052–2105)**. For instance, in Table 22, WizardMath-SFT exceeds Xwin-Math by 4.4% (48.1% vs. 43.7%) when using the Mistral-7B-v0.1 as the base model.

- **WizardMath vs. Skywork-Math:**

WizardMath-SFT underperforms Skywork-Math-2500k on the MATH benchmark by 12.1%, but it uses only 16.7% of the amount of data used by Skywork-Math-2500k (418k vs. 2500k), which is much less than Skywork-Math. Furthermore, according to **Figure 5 About Synthetic Data Size in the Skywork-Math paper[1]**, Skywork-Math-720k scores 34.54% on MATH, and Skywork-Math-360k scores 29.36%. Therefore, WizardMath-SFT-418k performs comparably to Skywork-Math-720k on MATH, and with the same amount of data, WizardMath-SFT outperforms Skywork-Math.

In summary, the Math Evol Instruct data synthesis method proposed in our study is as effective and practical as the current state-of-the-art data synthesis methods, such as DART-Math, Skywork-Math and Xwin-Math in the SFT stage. It significantly enhances the mathematical reasoning capabilities of the model, marking a key contribution of our work. Additionally, we acknowledge the contributions of methods such as DART-Math, Skywork-Math and Xwin-Math, which are excellent data synthesis approaches excelling in generating high-quality datasets for mathematical tasks and significantly enhancing models’ mathematical reasoning capabilities.

[1] Zeng L, Zhong L, Zhao L, et al. Skywork-Math: Data Scaling Laws for Mathematical Reasoning in Large Language Models–The Story Goes On[J]. arXiv preprint arXiv:2407.08348, 2024.

**Q2: This suggests that the main contribution of the paper is in the RL component. Therefore, the primary focus should be on the results obtained with different reward models, as presented in Table 4, utilizing various SFT backbones.**

Thank you for your deep insights. The following table 25 shows the impact of applying the proposed Instruction Quality Scoring Reward Model (IRM) and Process Supervised Reward Model (PRM) to PPO training across various SFT backbones (i.e., DART-Math, MetaMath, and Xwin-Math). The results demonstrate that incorporating our IRM and PRM during PPO training led to a performance improvement of 5% to 8% on both GSM8k and MATH for most SFT models. For instance:

- **When using DART-Math as the SFT backbone based on Llama2-7B:**

On GSM8k, after reinforcement learning training with IRM and PRM, Prop2Diff-RL improved by 6.9% (69.9% vs. 76.8%), and Uniform-RL improved by 5.3% (73.8% vs. 79.1%).

On MATH, Prop2Diff-RL achieved a 6.4% gain (30.7% vs. 37.1%), and Uniform-RL improved by 5.7% (29.5% vs. 35.2%).

- **When using DART-Math as the SFT backbone based on Mistral-7B-v0.1:**

On GSM8k, Prop2Diff-RL improved by 6.4% (81.1% vs. 87.5%), and Uniform-RL increased by 5.5% (82.6% vs. 88.1%).

On MATH, Prop2Diff-RL rose by 5.9% (45.5% vs. 51.4%), and Uniform-RL saw a 5.2% enhancement (43.5% vs. 48.9%).

- **For the MetaMath models based on Llama2-7B and Mistral-7B-v0.1:**

Training with PPO using IRM and PRM led to performance improvements of 8% to 9% on GSM8k and 5% to 8% on MATH.

- **Similarly, for the Xwin-Math-Llama2-7B model, performance on both GSM8k and MATH improved by 6% to 8%.**

These findings highlight the significant contributions of our IRM and PRM during reinforcement learning, consistently enhancing mathematical reasoning abilities of our SFT models while achieving robust generalization on different SFT backbones. This represents a key contribution of our study.

Thus, our study primarily makes two core contributions:

1. The proposed Math Evol Instruct data synthesis method is also as effective and practical as the current state-of-the-art data synthesis methods, such as DART-Math, Skywork-Math and Xwin-Math in the SFT stage. It also significantly enhances the mathematical reasoning capabilities of our models.

2. The proposed IRM and PRM models substantially improve performance during the reinforcement learning phase. They not only continuously enhance the mathematical reasoning abilities of our SFT models but also achieve strong generalization across various SFT backbones (i.e., DART-Math). Outstanding performance is demonstrated on the GSM8k and MATH.

**We have added the discussions about the Weaknesses-1.3 in Appendix C.1.4 of our latest upload of revised paper (pages 41–43, lines 2187–2309)**

Table 24: The performance comparison of WizardMath-SFT with DART-Math, Xwin-Math, and Skywork-Math on the Llama2-7B base model on the MATH benchmark.

Llama2 7B as the base model	Data size	MATH
DART-Math-Uniform	591k	29.5
DART-Math-Prop2Diff	585k	30.7
Xwin-Math	1440k	40.6
Skywork-Math	360k	29.36
Skywork-Math	720k	34.54
Skywork-Math	2500k	47.7
WizardMath-SFT	418k	35.6

Table 25: The impact of applying the proposed Instruction Quality Scoring Reward Model (IRM) and Process Supervised Reward Model (PRM) to PPO training across various SFT backbones (i.e., DART-Math, MetaMath, and Xwin-Math)

Model	Base	Params	GSM8k	MATH
MetaMath-SFT	Llama-2	7B	66.5	19.8
MetaMath-RL	Llama-2	7B	75.6	25.1
DART-Math-Prop2Diff-SFT	Llama-2	7B	69.9	30.7
DART-Math-Prop2Diff-RL	Llama-2	7B	76.8	37.1
DART-Math-Uniform-SFT	Llama-2	7B	73.8	29.5
DART-Math-Uniform-RL	Llama-2	7B	79.1	35.2
Xwin-Math-SFT	Llama-2	7B	82.6	40.6
Xwin-Math-RL	Llama-2	7B	88.2	48.5
WizardMath-Llama-SFT	Llama-2	7B	77.4	35.6
WizardMath-Llama-RL	Llama-2	7B	84.1	43.5
MetaMath-SFT	Mistral-v0.1	7B	77.9	28.6
MetaMath-RL	Mistral-v0.1	7B	86.4	35.2
DART-Math-Prop2Diff-SFT	Mistral-v0.1	7B	81.1	45.5
DART-Math-Prop2Diff-RL	Mistral-v0.1	7B	87.5	51.4
DART-Math-Uniform-SFT	Mistral-v0.1	7B	82.6	43.5
DART-Math-Uniform-RL	Mistral-v0.1	7B	88.1	48.7
WizardMath-Mistral-SFT	Mistral-v0.1	7B	82.8	48.1
WizardMath-Mistral-RL	Mistral-v0.1	7B	90.7	55.4

### C.1.5 WEAKNESSES-1.4

**Table 7 lacks adequate baselines; at least, the authors should include LLaMA-2-7B trained on the DartMath training set. This table also suffers from the same fairness issues as Table 1.**

Thank you for your constructive feedback. The table 26 below presents the performance of WizardMath-SFT on 7 out-of-domain (OOD) evaluation tasks covering K-12, college, and competition-level math problems in the SFT stage. The results indicate that **WizardMath-SFT** consistently surpasses state-of-the-art open-source models (i.e., DART-Math, Xwin-Math, and MathScale) across various scales and tasks, achieving an average improvement of **3%-6%**. For instance:

- With the Llama2-7B base model, WizardMath-SFT outperformed DART-Math-Uniform by **11.0%** (38.3% vs. 27.3%) and DART-Math-Prop2Diff by **10.5%** (38.3% vs. 27.8%) on average.

Table 26: The performance of WizardMath-SFT on the 7 out-of-domain evaluation results covering K-12, college, and competition level math problems compared with some SOTA models (i.e., DART-Math) in the SFT stage. The results of models in the table refer to MWPBENCH (Tang et al., 2024). “AGIE” stands for AGIEval. We report the models’ CoT pass@1 results on MwpBench without using any external python tool

Models	College Math	TAL	Math23k	Ape210k	Gaokao Bench Math	AGIE Gaokao Math	AGIE SAT Math	AVG
<i>Proprietary models</i>								
GPT-4	<b>24.4</b>	<b>51.8</b>	<b>76.5</b>	<b>61.5</b>	<b>35.4</b>	<b>28.2</b>	<b>68.6</b>	<b>49.5</b>
GPT-3.5-Turbo	21.6	42.9	62.5	44.0	23.2	15.3	55.8	37.9
<i>Models based on LLaMA-2 13B</i>								
LLaMA-2 13B	1.2	6.3	9.5	7.9	0.7	0.4	6.8	4.7
MAmmoTH-CoT	6.5	17.3	39.5	28.1	5.9	4.9	20.5	17.5
GAIR-Abel	7.9	21.1	42.2	27.8	7.0	4.9	30.3	20.2
MetaMath	10.1	25.4	48.6	31.6	9.6	5.6	38.2	24.2
MathScale 13B	20.4	38.1	61.1	43.7	20.0	12.3	55.8	35.9
WizardMath-SFT	22.2	42.5	65.9	47.6	31.6	23.5	59.7	41.9
WizardMath-RL	<b>22.9</b>	<b>43.3</b>	<b>70.3</b>	<b>50.8</b>	<b>33.1</b>	<b>25.7</b>	<b>64.7</b>	<b>44.4</b>
<i>Models based on LLaMA-2 7B</i>								
LLaMA-2 7B	2.3	7.6	6.8	7.3	2.1	2.9	2.9	4.6
MAmmoTH-CoT	6.2	13.3	34.6	21.4	3.9	2.7	19.6	14.5
GAIR-Abel	6.6	18.3	35.4	24.5	4.3	4.4	23.5	16.7
MetaMath	9.4	22.5	44.0	29.9	5.9	5.1	36.2	21.9
DART-Math-Uniform	12	27.3	47.9	32.9	14.8	11.1	45.1	27.3
DART-Math-Prop2Diff	11.9	27.7	49.9	34.3	12.8	10.6	47.1	27.8
Xwin-Math-V1.1	14.9	29.7	59.6	40.8	15.9	8.4	51.0	31.5
MathScale 7B	20.9	35.2	59.0	41.8	19.6	12.6	57.8	35.3
WizardMath-SFT	21.1	38.5	62.4	43.8	26.3	17.7	58.3	38.3
WizardMath-RL	<b>21.2</b>	<b>40.2</b>	<b>67.3</b>	<b>46.1</b>	<b>28.9</b>	<b>18.7</b>	<b>62.7</b>	<b>40.7</b>
<i>Models based on Mistral 7B</i>								
Mistral 7B	7.5	17.9	18.5	15.5	6.2	5.9	22.5	13.4
MetaMath Mistral	15.7	31.4	55.1	38.1	15.3	10.1	50.9	30.9
DART-Math-Uniform	19.4	34.8	61.6	44.8	27.0	16.1	59.8	37.6
MathScale Mistral	21.8	39.9	64.4	46.0	21.4	14.3	57.8	37.9
DART-Math-Prop2Diff	19.9	37.4	62.2	44.9	27.2	18.1	62.7	38.9
WizardMath-Mistral-SFT	24.3	42.7	66.6	49.7	35.2	22.7	63.1	43.5
WizardMath-Mistral-RL	<b>24.8</b>	<b>44.8</b>	<b>71.2</b>	<b>52.6</b>	<b>37.2</b>	<b>24.5</b>	<b>64.7</b>	<b>45.7</b>

- With the Mistral-7B base model, WizardMath-SFT achieved an average improvement of **5.9%** over DART-Math-Uniform (43.5% vs. 37.6%) and **4.6%** over DART-Math-Prop2Diff (43.5% vs. 38.9%).

These findings highlight the effectiveness of our **Math Evol-Instruct** method, demonstrating its robustness and superior generalization capabilities on out-of-domain tasks.

**We have added the discussions about the Weaknesses-1.4 in Appendix C.1.5 of our latest upload of revised paper (pages 43–44, lines 2310–2363)**

## C.2 RECOMMEND

**I recommend that the authors reorganize the paper better to emphasize their contributions to the "RL part."**

Thank you for your deep insights and constructive suggestions. Due to time and space constraints, we promise to further emphasize our contributions to the "RL part" in future revisions of our paper. Specifically, we will provide more detailed descriptions of the contributions of our proposed RLEIF approach to the RL part in some sections (i.e., the Abstract, Introduction, and Experiment Sections). For instance, we will highlight that in RL training, we firstly propose the instruction quality scoring reward model combined with the process supervision reward model not only continuously enhancing the mathematical reasoning abilities of the SFT model but also achieve strong generalization across various SFT backbones. Additionally, we will supplement the discussion on the application and

2376 impact of IRM and PRM on different advanced SFT backbones, as highlighted in the Weaknesses-1.3,  
2377 to further strengthen the theoretical framework and experimental analysis.

### 2378 2379 C.3 QUESTIONS 2380

2381 Thank you very much for your insightful questions and valuable suggestions. Below, we provide  
2382 responses to your Question-1 through Question-5 in sequence.

#### 2383 2384 C.3.1 QUESTIONS-1

2385 **In Equation (1), how is the parameter ( $m$ ) set? Additionally, how is the Instruction Reward  
2386 Model (IRM) trained?**

2387  
2388 **Q1: In Equation (1), how is the parameter ( $m$ ) set?**

2389 The parameter ( $m$ ) denotes the margin in the Pairwise Ranking Loss, acting as a threshold to regulate  
2390 the score difference between <Choose, Reject> pairs. Specifically, it ensures that during IRM training,  
2391 the reward score for higher-quality instructions surpasses that of lower-quality instructions by at least  
2392 the margin value. This mechanism encourages the model to emphasize the quality score gap between  
2393 high-quality and low-quality instructions. In our experiments, the parameter ( $m$ ) was set to a constant  
2394 1.

2395 **Q2: Additionally, how is the Instruction Reward Model (IRM) trained?**

2396 **In our paper, Section 3.2 <REWARD MODELS>, lines 187-201,** we conducted two rounds  
2397 of downward evolution and three rounds of upward evolution based on the original instructions,  
2398 generating a total of five evolved instructions. Subsequently, we leverage GPT-4 to rank the quality  
2399 between those evolved instructions and original instruction based on the difficulty and definition,  
2400 with higher ranks assigned to instructions demonstrating greater difficulty and clearer definitions.  
2401 The detailed ranking prompt template is provided in Appendix A.2.

2402 From the ranking results of the 6 instructions, we created 15 positive-negative sample pairs by  
2403 combining  $C(6, 2)$ . Applying this five-round evolution process to 15k original instructions, we  
2404 ultimately generated  $15k \times 15 = 225k$  positive-negative pairs for training IRM data.

2405 During training, we employed the Pairwise Ranking Loss defined in Eq. 1. For a given mathematical  
2406 instruction  $g$ , the IRM quantifies its quality by assigning a score. The IRM was initialized with  
2407 the SFT model and augmented with a header layer that outputs a scalar score. The design of the  
2408 Pairwise Ranking Loss draws inspiration from the reward model training methods described in the  
2409 Instruct-GPT paper[1].  
2410

2411 **We have added the answers about the Questions-1 in Appendix C.3.1 of our latest upload of  
2412 revised paper (page 45, lines 2383–2415)**

2413 [1] Ouyang L, Wu J, Jiang X, et al. Training language models to follow instructions with human  
2414 feedback[J]. Advances in neural information processing systems, 2022, 35: 27730-27744.

#### 2415 2416 C.3.2 QUESTIONS-2

2417 **Lines 256-258 suggest that you retain solutions with incorrect answers. How might this influence  
2418 the results? Have you considered using the IRM to filter out low-quality examples for supervised  
2419 fine-tuning (SFT)?**

2420  
2421 **Q1: Lines 256-258 suggest that you retain solutions with incorrect answers. How might this  
2422 influence the results?**

2423 **The Lines 256-258, Sections 4.2 for SFT Training Data in our paper** In order to prevent data  
2424 leakage, we filter out the evolved data with high similarity to the GSM8k and MATH test sets,  
2425 so it does not refer to incorrect answers. The data leakage detection method refers to the paper  
2426 Appendix A.11 line 1782-1815. Specifically, we employ instructions in the GSM8k and MATH test  
2427 set as queries to retrieve the top-5 samples from all evolved training data with an embedding model,  
2428 gte-large. Additionally, we employ GPT-4 to provide similarity judgement between the test sets and  
2429 the retrieved samples, and remove the similar instructions. The prompt and additional details are  
provided in Appendix A.12.

2430 The table 27 below demonstrates the impact of unfilter the potential data leaks on model performance.  
 2431 WizardMath-SFT-No-filter-data-leakage outperforms WizardMath-SFT-Filter-data-leakage by 1.3%  
 2432 on the GSM8k and by 1.7% on the MATH. we use Mistral-7B-v0.1 as the base model  
 2433

2434  
 2435 Table 27: The impact of unfilter the potential data leaks on model performance. we use Mistral-7B-  
 2436 v0.1 as the base model

Model	GSM8k	MATH
WizardMath-SFT-Filter-data-leakage	82.8	48.1
WizardMath-SFT-No-Filter-data-leakage	84.1	49.8

2440  
 2441 **Q2: Have you considered using the IRM to filter out low-quality examples for supervised**  
 2442 **fine-tuning (SFT)?**

2443 Thank you very much for your insightful question and constructive suggestions. The table 28 below  
 2444 highlights the effects of using IRM to filter out low-quality instructions during the SFT stage.  
 2445

- 2446 • Filtering 15k low-quality instructions resulted in **WizardMath-SFT-filter-15k** outperforming  
 2447 WizardMath-SFT-original, with a 1.8-point improvement on GSM8k and a 2.1-point  
 2448 improvement on MATH.
- 2449 • Filtering 30k low-quality instructions improved GSM8k by 0.9% and MATH by 0.6%.
- 2450 • However, when the filtering reached 45k, WizardMath-SFT-filter-45k showed a performance  
 2451 decrease of 0.8% on GSM8k and 1.1% on MATH.
- 2452 • Filtering up to 60k resulted in a more pronounced decline, with WizardMath-SFT-filter-60k  
 2453 dropping by 1.7% on GSM8k and 2.5% on MATH.  
 2454

2455 These results indicate that using IRM for moderate filtering of low-quality data (i.e., 15k or 30k)  
 2456 is effective for enhancing model performance, while excessive filtering can lead to significant  
 2457 performance degradation.

2458 **We have added the discussions about the Questions-2 in Appendix C.3.2 of our latest upload of**  
 2459 **revised paper (pages 45-46, lines 2416–2470)**  
 2460

2461  
 2462 Table 28: The impact of employing IRM to filter low-quality data on model performance in the SFT  
 2463 stage. we use Mistral-7B-v0.1 as the base model

Model	IRM Filter Data Size	GSM8k	MATH
WizardMath-SFT-original	-	82.8	48.1
WizardMath-SFT-filter-15k	15k	84.6	50.2
WizardMath-SFT-filter-30k	30k	83.7	48.7
WizardMath-SFT-filter-45k	45k	82.0	47.0
WizardMath-SFT-filter-60k	60k	81.1	45.6

### 2471 C.3.3 QUESTIONS-3

2472  
 2473 **During PPO, do you use two reward models? Using two reward models in PPO can be time-**  
 2474 **consuming and computationally expensive. What are your strategies for addressing this?**  
 2475

2476 **Q1: During PPO, do you use two reward models?**

2477 Yes, we use two reward models during the PPO training stage.

2478  
 2479 **Q2: Using two reward models in PPO can be time-consuming and computationally expensive.**  
 2480 **What are your strategies for addressing this?**

2481 During the PPO training stage, we utilized the DeepSpeedChat[1] framework. To improve training  
 2482 efficiency and reduce memory consumption, we employed several optimization techniques,  
 2483 including DeepSpeed ZeRO-3 with CPU Offload, the DeepSpeed-Hybrid Engine, MixZ++, Gradient  
 Checkpointing, Gradient Accumulation, and BFloat16 precision.

2484 In the future we try to implement GRPO by Deepseekmath[2] (a variant of PPO) for training and  
 2485 incorporate the VLLM[3] used by the OpenRLHF[4] framework to accelerate the policy model  
 2486 generation during PPO training, thus improving the training efficiency.

2487 **We have added the answer about the Questions-3 in Appendix C.3.3 of our latest upload of**  
 2488 **revised paper (pages 46-47, lines 2472–2497)**

2490 [1] Yao Z, Aminabadi R Y, Ruwase O, et al. Deepspeed-chat: Easy, fast and affordable rlhf training  
 2491 of chatgpt-like models at all scales[J]. arXiv preprint arXiv:2308.01320, 2023.

2492 [2] Shao Z, Wang P, Zhu Q, et al. Deepseekmath: Pushing the limits of mathematical reasoning in  
 2493 open language models[J]. arXiv preprint arXiv:2402.03300, 2024.

2494 [3] Kwon W, Li Z, Zhuang S, et al. Efficient memory management for large language model serving  
 2495 with pagedattention[C]//Proceedings of the 29th Symposium on Operating Systems Principles. 2023:  
 2496 611-626.

2498 [4] Hu J, Wu X, Wang W, et al. OpenRLHF: An Easy-to-use, Scalable and High-performance RLHF  
 2499 Framework[J]. arXiv preprint arXiv:2405.11143, 2024.

#### 2500 C.3.4 QUESTIONS-4

2502 **In Lines 89-90, you state that you "innovatively introduce PRM to address the False-Positive**  
 2503 **issue in the problem-solving process." This claim should be validated by comparing the false-**  
 2504 **positive rate on a test set both with and without your method.**

2505 Thank you for your insightful feedback. We utilized GPT-4o-2024-0513 (the accuracy is 96.1% On  
 2506 GSM8k and 76.6% on MATH) to annotate the step-by-step correctness of responses generated by  
 2507 WizardMath-SFT, WizardMath-RL-ORM, and WizardMath-RL-PRM on the GSM8k and MATH test  
 2508 sets, and we calculated the model's false-positive rates.

2509 We define the **false-positive rate** as the proportion of responses in the test set where the final answer  
 2510 is correct, but errors occur in intermediate steps (i.e., computational or logical mistakes). The formula  
 2511 for calculating the False Positive Rate is as follows:

$$2512 \text{False Positive Rate} = \frac{\text{Number of False Positives}}{\text{Total Number of Test Sets}}$$

2516 The table 29 below presents the statistical results:

- 2518 • The false-positive rate of WizardMath-SFT is 2.58% on GSM8k and 2.36% on MATH.
- 2519 • The false-positive rate of WizardMath-RL-ORM is 1.67% on GSM8k and 1.56% on MATH.
- 2520 • The false-positive rate of WizardMath-RL-PRM is 0.68% on GSM8k and 0.90% on MATH.

2522 Compared to WizardMath-SFT, WizardMath-RL-PRM reduced the false-positive rate by 1.90% on  
 2523 GSM8k and 1.46% on MATH. Similarly, compared to WizardMath-RL-ORM, WizardMath-RL-PRM  
 2524 achieved reductions of 0.99% on GSM8k and 0.66% on MATH.

2526 These results demonstrate that the incorporation of PRM significantly reduces the model's false-  
 2527 positive rate, effectively alleviating the occurrence of intermediate step errors during the problem-  
 2528 solving process.

2529 **We have added the discussions about the Questions-4 in Appendix C.3.4 of our latest upload of**  
 2530 **revised paper (page 47, lines 2500–2528)**

#### 2532 C.3.5 QUESTIONS-5

2533 **In Lines 88-89, you mention that existing methods "mainly focus on the SFT stage and are**  
 2534 **susceptible to learning hallucinated information from the teacher model." However, in Line 95,**  
 2535 **you still use GPT-4 to annotate step-level labels. Isn't there a risk of obtaining incorrect step**  
 2536 **labels from GPT-4 as well?**

2537 Yes, there is a potential risk of obtaining incorrect step labels from GPT-4 as well

Table 29: The impact of PRM to alleviate the False-Positive issue in the RL training stage. we use Mistral-7B-v0.1 as the base model

Metrics	WizardMath-SFT	WizardMath-RL-ORM	WizardMath-RL-PRM
Reward Model for PPO	-	ORM	PRM
Number of GSM8k test sets	1319	1319	1319
Number of False Positive On GSM8k	34	22	9
False Positive Rate On GSM8k	2.58%	1.67%	0.68%
Number of MATH test sets	5000	5000	5000
Number of False Positive On MATH	118	78	45
False Positive Rate On MATH	2.36%	1.56%	0.90%

The risk of the model learning hallucinatory information from the teacher model cannot be completely eliminated. Therefore, in order to ensure the reliability and effectiveness of using GPT-4 to annotate step-level labels during the problem-solving process in constructing PRM training data, we conducted the following two analyses:

### 1. Reliability of GPT-4 Annotations

Manually annotating large-scale step-level PRM training data demands extensive mathematical expertise, making it a challenging, time-intensive, and costly process. So, we employ a fully AI-powered automatic annotation using GPT-4 in our paper. To assess the reliability of GPT-4-generated annotations, in the early stages, we randomly selected 2k samples from the manually labeled PRM800k step-level training dataset and annotated them using GPT-4. GPT-4 annotations were evaluated against human annotations using the F1 score as a consistency metric. The results showed an F1 consistency of 78.1% between GPT-4 and human annotations.

Additionally, for the GSM8k training set, which is relatively lower in difficulty, we randomly sampled 200 examples for step-level labeling using GPT-4 and manual annotations. The results show that the F1 consistency between GPT-4 and manual labeling on GSM8k is 87.2%. These findings demonstrate that the annotation using GPT-4 with manual evaluation exhibits high consistency on GSM8k and MATH, thus ensuring the reliability of step-level annotation using GPT-4 for PRM training data.

### 2. Effectiveness of GPT-4 Annotations

The table 30 below and Table 4 in our paper (lines 382–395, 415–424) discussed the impact of AI-labeled PRM data on model performance compared to manually labeled PRM800k and Math-Shepherd data generated via MCTS Tree Search. The experimental results reveal that the PRM trained on our fully AI-labeled data outperforms both the manually annotated PRM800k and Math-Shepherd. For instance:

- When training WizardMath-Llama2-7B-SFT with PPO, GPT-4-labeled PRM data surpasses PRM800k by 2.0% and Math-Shepherd by 1.4% on GSM8k, and by 1.2% and 1.7%, respectively, on MATH.
- Similarly, with WizardMath-Mistral-7B-SFT trained using PPO, GPT-4-labeled PRM data outperforms PRM800k by 1.8% and Math-Shepherd by 1.1% on GSM8k, and by 1.9% and 2.4%, respectively, on MATH.

Moreover, PRM outperforms ORM by 2%–3% on both GSM8k and MATH, achieving a notable improvement of 4%–5% on WizardMath-SFT. These results highlight the effectiveness of GPT-4-labeled data for PRM training. (It is worth noting that our evolved instructions lack correct answers, limiting compatibility with the methods employed by Math-Shepherd which needs the correct answers.)

The analysis above underscores both the reliability and effectiveness of using GPT-4 to annotate step-level PRM training data. However, we acknowledge that GPT-4 annotations are not immune to errors, and the possibility of incorrect step labels represents a limitation of this approach.

**We have added the discussions about the Questions-5 in Appendix C.3.5 of our latest upload of revised paper (pages 47-49, lines 2532–2590)**



Table 30: The effect of using manually labeled and AI labeled PRM training data in PPO training

Models	GSM8K	MATH
Llama2-7B: WizardMath-SFT	77.4	35.6
+ ORM (ours)	79.1	36.8
+ PRM800k	79.7	38.7
+ Math-Shepherd	80.3	38.2
+ PRM (ours)	<b>81.7</b>	<b>39.9</b>
Mistral-7B: WizardMath-SFT	82.8	48.1
+ ORM (ours)	84.6	49.6
+ PRM800k	85.4	50.8
+ Math-Shepherd	86.1	50.3
+ PRM (ours)	<b>87.2</b>	<b>52.7</b>

#### C.4 REVIEWER-DHFE

Dear Reviewer dHFe,

We sincerely thank you for your insightful comments and the time you dedicated to reviewing our work. Your expert feedback has been invaluable in guiding us towards refining our paper and making it more comprehensive and competitive. We greatly appreciate your support and constructive suggestions. In the following, we offer detailed responses to the **Weaknesses** and **Questions** raised in your review, addressing each point in a systematic manner.

Furthermore, in **Appendix C.4 of our latest upload of revised paper (pages 49–57, lines 2606–3077)**, we also have added the discussions with the Reviewer-dHFe on the weaknesses and questions of our paper to respond to the Reviewer-dHFe’s comments and to further improve the quality of our research.

##### C.4.1 WEAKNESSES-1

**PRM labels from GPT-4 – Not really sure what to think of this. On one hand, I feel such direct distillation like this would limit the effectiveness of a method at larger data scales. On the other hand, the results seem to be good (and also this is one key part that makes the process fully AI-automated.)**

Thank you very much for your constructive feedback and recognition of the effectiveness of our approach. To explore the effectiveness and reliability of using GPT-4 to annotate PRM training data at larger data scales, we conducted a Data Scaling law analysis of using GPT-4 to annotate PRM training data. Additionally, we explore the feasibility of leveraging open-source models (i.e., Llama-3.1-405B-Instruct) as cost-effective alternatives to GPT-4 for annotation.

#### 1. Effectiveness Analysis of the Scaling Law in GPT-4 Annotated PRM Training Data.

##### 1.1 Impact of PRM Data Scaling when PRM as the Verifier for the Best of N Metric.

To assess the influence of data scale when PRM acts as the verifier, we randomly sampled subsets of 50k, 150k, and 300k data from our total 450k PRM training dataset as shown in the following table 31. Models were trained on these subsets, and we evaluated the Best of N metric on GSM8k and MATH. Following the same settings as Table 5 in our paper, we sampled 256 answers for each problem, scored them with the PRM Verifier, and selected the highest-scoring answer. Key results are summarized as follows:

- On GSM8k, the Best of N performance of PRM significantly improved as the training data size increased. For instance, PRM-450k achieved 95.2%, outperforming PRM-300k by 1.6% and PRM-150k by 2.9%.
- On MATH, PRM-450k reached 64.7%, marking a 1.4% improvement over PRM-300k and a 3.2% improvement over PRM-150k.

Table 31: Impact of PRM Data Scaling when PRM as the Verifier for the Best of N Metric. The SFT and Reward models are trained based on Mistral-7B. The verifier is based on 256 sample outputs.

Generators	Verifiers	PRM Data Size	GSM8K	MATH
WizardMath-SFT	PRM-50k	50k	89.6	58.9
	PRM-150k	150k	92.3	61.5
	PRM-300k	300k	93.6	63.3
	PRM-450k	450k	95.2	64.7

## 1.2 Impact of PRM Data Scaling on PPO Training Performance.

In the following table 32, we further investigated the effect of PRM training data scaling during the PPO training stage. Increasing the PRM data size yielded substantial performance gains for WizardMath-RL on GSM8k and MATH:

- On GSM8k, PPO training with PRM-450k achieved 87.2%, surpassing PRM-300k by 1.4% and PRM-50k by 3.7%.
- On MATH, PPO training with PRM-450k reached 52.7%, exceeding PRM-300k by 1.5% and PRM-50k by 4.0%.

These findings confirm that scaling PRM training data consistently enhances PRM performance as a Verifier on the Best of N metric and significantly improves PPO training outcomes. This validates the effectiveness of GPT-4 annotated PRM training data in adhering to the Data Scaling Law, demonstrating their robustness and utility even at larger data scales.

Table 32: The effect of PRM training data scaling during the PPO training stage. We use Mistral-7B-v0.1 as the base model

Models	GSM8K	MATH
Mistral-7B: WizardMath-SFT	82.8	48.1
+ PRM-50k	83.5	48.7
+ PRM-150k	84.9	49.8
+ PRM-300k	85.8	51.2
+ PRM-450k	<b>87.2</b>	<b>52.7</b>

## 2. Reliability analysis of GPT-4 labeled PRM training data compared to manual labeling

To assess the reliability of GPT-4-generated annotations, in the early stages, we randomly selected 2k samples from the manually labeled PRM800k step-level training dataset and annotated them using GPT-4. GPT-4 annotations were evaluated against human annotations using the F1 score as a consistency metric. The results showed an F1 consistency of 78.1% between GPT-4 and human annotations.

Additionally, for the GSM8k training set, which is relatively lower in difficulty, we randomly sampled 200 examples for step-level labeling using GPT-4 and manual annotations. The results show that the F1 consistency between GPT-4 and manual labeling on GSM8k is 87.2%. These findings demonstrate that the annotation using GPT-4 with manual evaluation exhibits high consistency on GSM8k and MATH, thus ensuring the reliability of step-level annotation using GPT-4 for PRM training data.

## 3. Feasibility of using advanced open-source models instead of GPT-4 to label PRM training data

We realize that there is a high cost of directly distilling GPT-4 in large-scale data scenarios, which is a limitation of this study. Additionally, manual annotation demands mathematical expertise and entails a challenging, time-intensive, and costly process. Moreover, our evolved instructions lack correct answers, limiting compatibility with the methods employed by Math-Shepherd[1] which needs the correct answers.

To mitigate these challenges, we also explore the feasibility of leveraging advanced open-source models, such as Llama-3.1-405B-Instruct, instead of GPT-4 for PRM training data labeling, using the

Table 33: The impact of using advanced open-source models(i.e., Llama-3.1-405B-Instruct) for PRM training data labeling and we use Mistral-7B-v0.1 as the base model.

Models	AI-Label	GSM8k	MATH
WizardMath-SFT	-	82.8	48.1
+ PRM-Llama-3.1-405B-Instruct	Llama-3.1-405B-Instruct	85.8	51.5
+ PRM-GPT-4	GPT-4	87.2	52.7

same label prompts and training settings. As shown in the table 33, WizardMath-PRM-Llama-3.1-405B achieves **85.8%** on the GSM8k, outperforming WizardMath-SFT by **3.0%** and lagging behind WizardMath-PRM-GPT-4 by **1.4%**. On the MATH, it scores **51.5%**, exceeding WizardMath-SFT by **3.4%** with a **1.2%** gap compared to WizardMath-PRM-GPT-4.

Balancing cost and accuracy, Llama-3.1-405B-Instruct demonstrates considerable potential as a substitute for GPT-4 in PRM training data labeling.

In conclusion, GPT-4-based labeled PRM data also follows the data scaling law and offers an effective solution for larger data scales. For scenarios requiring a balance between cost and accuracy, advanced open-source models like **Llama-3.1-405B-Instruct** provide a viable alternative. We hope these analyses above can address your concerns.

[1]. Wang P, Li L, Shao Z, et al. Math-shepherd: Verify and reinforce llms step-by-step without human annotations[C]//Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2024: 9426-9439.

#### C.4.2 WEAKNESSES-2

**Unclear presentation – The paper assumes that readers are already previously familiar with Evol-Instruct, as it devotes very little time to talking about it in the intro or related work. The narrative is messy – there are certain concepts (e.g. "grade school" and "high school" questions) that were introduced once out of nowhere then never mentioned again. There are a number of rows on Table 1 that were never discussed in the main text. Figure 1 is difficult to understand. There are also several typos and poorly-worded sentences.**

**Q1: Unclear presentation – The paper assumes that readers are already previously familiar with Evol-Instruct, as it devotes very little time to talking about it in the intro or related work.**

Thank you for highlighting the valuable questions and we sincerely apologize for any inconvenience caused. We have added an introduction to Evol-Instruct in our latest upload of revised paper (Section1 <INTRODUCTION>, Lines 50–82). A more comprehensive description can be found in Appendix C.4.2, with the relevant details outlined below. In future camera-ready version of our paper, we promise to integrate this section in the "Introduction" or "Related Work" sections and provide a more comprehensive explanation of Evol-Instruct. The relevant detailed descriptions as follows:

Evol-Instruct proposed by WizardLM[1] is an innovative framework designed to automate the generation of diverse and complex open-domain instructions using large language models (LLMs). Instead of relying on human-crafted instruction datasets, it leverages the generative capabilities of LLMs to iteratively evolve an initial set of instructions through two complementary strategies: In-depth Evolving and In-breadth Evolving. Starting from an initial dataset  $D^{(0)} = \{(I_k^{(0)}, R_k^{(0)})\}_{k=1}^N$ , it iteratively evolves instructions over  $M$  turns, producing datasets  $[D^{(1)}, \dots, D^{(M)}]$ .

Two evolution strategies are employed: In-depth Evolving incrementally enhances instruction complexity by introducing additional constraints, deepening, concretizing, increasing reasoning steps, and complicating input, while maintaining logical coherence and ensuring instructions remain solvable by humans,  $I_k^{(t+1)} = \text{In-Depth Operation}(I_k^{(t)})$ ; In-breadth Evolving focuses on improving topic diversity and dataset richness by creating entirely new instructions, expanding the coverage of skills and scenarios, particularly in underrepresented areas,  $I_k^{(t+1)} = \text{In-Breadth Operation}(I_k^{(t)})$ .

To ensure dataset quality, failed evolutionary instructions are filtered. Evol-Instruct supports scalable, high-quality dataset creation, significantly enhancing LLM performance in reasoning and open-domain tasks. Notably, WizardCoder[2] incorporates instruction evolution specifically tailored to coding tasks, leading to substantial improvements in code generation capabilities.

[1]. Xu C, Sun Q, Zheng K, et al. WizardLM: Empowering large pre-trained language models to follow complex instructions[C]//The Twelfth International Conference on Learning Representations. 2024.

[2]. Luo Z, Xu C, Zhao P, et al. Wizardcoder: Empowering code large language models with evol-instruct[J]. arXiv preprint arXiv:2306.08568, 2023.

**Q2: The narrative is messy – there are certain concepts (e.g. "grade school" and "high school" questions) that were introduced once out of nowhere then never mentioned again.**

We greatly appreciate you highlighting this writing issue and sincerely apologize for any inconvenience it may have caused. In our paper, we primarily evaluate the model’s mathematical performance on two popular benchmarks: GSM8k and MATH. GSM8k represents problems at the grade school level, while MATH focuses on high school competition problems, such as AMC 10, AMC 12, and AIME. In Appendix C.4.2, we have included a detailed introduction to the GSM8k and MATH datasets. Additionally, in our latest upload of revised paper, we have incorporated descriptions of "grade school" and "high school," as reflected in lines 90-91, 102-103, 114, 240, 340, and 537.

We also include a detailed description of the evaluation benchmark as follows:

We mainly evaluate WizardMath on two benchmarks (GSM8k and MATH). The GSM8k dataset contains approximately 7500 training data and 1319 test data, mainly on grade school level math problems, each of which consists of basic arithmetic operations (addition, subtraction, multiplication, and division), and generally requires 2 to 8 steps to solve. The MATH dataset collects math problems from prestigious math competitions such as AMC 10, AMC 12, and AIME. It contains 7500 training data and 5,000 challenging test data in seven academic areas: Prealgebra, Algebra, Number Theory, Counting and Probability, Geometry, Intermediate Algebra, and Precalculus. Furthermore, these problems are divided into five levels of difficulty, with ‘1’ denoting the relatively lower difficulty level and ‘5’ indicating the highest level.

**Q3: There are a number of rows on Table 1 that were never discussed in the main text.**

Thank you for pointing out this valuable questions. Below, we present a detailed analysis of the performance improvements across various model scales (0.1B to 70B) with different base models on GSM8k and MATH benchmarks:

#### (1.) Using GPT-2 Series Models as the Base Model:

- **GPT-2-Small-0.1B:** WizardMath-GPT-2-Small improves by 19.5% (26.4 vs. 6.9) on GSM8k and 6.9% (12.3 vs. 5.4) on MATH.
- **GPT-2-Medium-0.3B:** WizardMath-GPT-2-Medium enhances by 27.5% (38.7 vs. 11.2) on GSM8k and by 9.4% (15.6 vs. 6.2) on MATH, outperforming Llama2-13B.
- **GPT-2-Large-0.7B:** WizardMath-GPT-2-Large increases by 36.5% (50.1 vs. 13.6) on GSM8k and by 14.8% (21.2 vs. 6.4) on MATH, surpassing Mistral-7B-v0.1.
- **GPT-2-XL-1.5B:** WizardMath-GPT-2-XL shows a 43.5% (58.9 vs. 15.4) gain on GSM8k and 18.5% (25.4 vs. 6.9) on MATH, exceeding MAMmoTH-CoT-Llama2-13B.

These results demonstrate that the RLEIF method significantly enhances the mathematical reasoning capabilities of GPT-2 series base models.

#### (2.) Using Llama Series Models as the Base Model:

- **Llama-3.2-1B:** WizardMath-Llama-3.2-1B improves by 18.9% (63.3 vs. 44.4) on GSM8k and by 2.9% (33.5 vs. 30.6) on MATH compared to Llama-3.2-1B-Instruct.
- **Llama-3.2-3B:** WizardMath-Llama-3.2-3B enhances GSM8k by 7.8% (85.5 vs. 77.7) and MATH by 1.9% (49.9 vs. 48.0).

- 2808 • **Llama-2-7B:** WizardMath-Llama-2-7B achieves a 69.5% improvement on GSM8k (84.1 vs. 14.6) and 41.0% on MATH (43.5 vs. 2.5), surpassing Xwin-Math-Llama-2-7B, MathScale-Llama-2-7B, and MetaMath-Llama-2-7B.
- 2809
- 2810
- 2811 • **Llama-3-8B:** WizardMath-Llama-3-8B attains 90.3% on GSM8k (1.7% higher than Jiuzhang3.0) and 58.8% on MATH (7.8% higher than Jiuzhang3.0), also outperforming Baichuan-3, GLM-4, Gemini-Pro, Claude2, and GPT-3.5-Turbo, and is comparable to GPT-4-0314.
- 2812
- 2813
- 2814
- 2815 • **Llama-2-13B:** WizardMath-Llama-2-13B improves GSM8k by 61.0% (89.7 vs. 28.7) and MATH by 46.7% (50.6 vs. 3.9), outperforming SOTA models such as Xwin-Math and KPMath-Plus.
- 2816
- 2817
- 2818
- 2819 • **Llama-2-70B:** WizardMath-Llama-2-70B enhances GSM8k by 36.0% (92.8 vs. 56.8) and MATH by 45.1% (58.6 vs. 13.5).
- 2820

### 2821 (3.) Using Mistral Series Models as the Base Model:

- 2822
- 2823 • **Mistral-7B-v0.1:** WizardMath-Mistral-7B-v0.1 improves GSM8k by 47.8% (90.7 vs. 42.9) and MATH by 42.5% (55.4 vs. 12.9).
- 2824
- 2825 • **Mistral-7B-v0.3:** WizardMath-Mistral-7B-v0.3 achieves 90.4% on GSM8k and 55.6% on MATH, comparable to WizardMath-Mistral-7B-v0.1.
- 2826
- 2827
- 2828 • **Mathstral-7B-v0.1:** WizardMath-Mathstral-7B-v0.1 attains 93.8% on GSM8k and 70.9% on MATH, comparable to GPT-4-Turbo-0125 and Claude 3.5 Sonnet, and superior to GPT-4 (original version).
- 2829
- 2830

2831 **(4.) Using DeepSeekMath as the Base Model:** WizardMath-DeepSeek improves GSM8k by 26.8% and MATH by 28.4%, outperforming DART-Math and DeepSeekMath-RL.

### 2832 (5.) Using Qwen2.5 Series Models as the Base Model:

- 2833
- 2834
- 2835 • **Qwen2.5-Math-2.5B:** WizardMath-Qwen2.5-Math-2.5B achieves 86.7% on GSM8k and 68.6% on MATH.
- 2836
- 2837 • **Qwen2.5-Math-7B:** WizardMath-Qwen2.5-Math-7B attains 93.9% on GSM8k and 77.8% on MATH.
- 2838
- 2839
- 2840 • **Qwen2.5-7B:** WizardMath-Qwen2.5-7B achieves 94.0% on GSM8k and 74.5% on MATH, performing comparably to GPT-4o-2024-0516 and Claude 3.5 Sonnet.
- 2841

2842 The proposed RLEIF method significantly enhances the mathematical reasoning performance across various scales ranging from 0.1B to 70B with different base models, consistently outperforming all state-of-the-art open-source models. Notably, WizardMath-Mathstral-7B-v0.1 and WizardMath-Qwen2.5-Math-7B surpass some proprietary models such as GPT-4 (original version), Gemini-Pro, and GPT-3.5-Turbo, and perform comparably to GPT-4-Turbo-0125, GPT-4o-2024-0516, and Claude 3.5 Sonnet. These findings further underscore the effectiveness, robustness, and scalability of the proposed RLEIF method in our study.

2849 In future camera-ready versions of our paper, we promise to incorporate the results presented above into Section 4.3 <Main Results> and Appendix C.4.2, to provide a more comprehensive and in-depth analysis of the effectiveness of the proposed RLEIF method in enhancing the model’s mathematical reasoning capabilities across a range of model scales (0.1B to 70B) with various base models.

### 2853 Q4: Figure 1 is difficult to understand.

2854

2855 Thank you for highlighting these important questions and for pointing out any confusion caused by Figure 1. We sincerely apologize for any lack of clarity and greatly appreciate your valuable feedback. In future camera-ready versions of our paper, we are committed to providing a more comprehensive explanation of the diagram. Furthermore, we offer a detailed clarification of our method flow below, including the significance of colors and shapes as well as the direction of the arrows in Figure 1, to facilitate clearer understanding.

2858

2859

2860

2861 In Figure 1, the various colored squares represent specific elements: **blue squares** denote original instructions, **orange squares** indicate evolved instructions, **cyan squares** signify model-generated

solution processes, and **grey squares** correspond to a series of training-related operations such as supervised fine-tuning (SFT), reward modeling, and reinforcement learning (RL). To enhance the mathematical reasoning capabilities of large language models, we propose the **RLEIF method**, which integrates instruction evolution with reinforcement learning. This method consists of three primary steps:

### 1. Instruction Evolution and SFT

In the first step, we apply upward and downward instruction evolution on the GSM8k and MATH datasets, generating evolved instructions for the SFT. On the leftmost side of Figure 1, the three blue arrows, from top to bottom, represent:

- (a) the adoption of the instruction evolution technique,
- (b) the generation of evolved instruction data, and
- (c) its application to SFT training.

### 2. Reward Model Training

The second step involves two reward models: the **Instruction Quality Scoring Reward Model (IRM)** and the **Process-Supervised Reward Model (PRM)**, depicted in the central section of Figure 1.

- **IRM:** We employ upward and downward evolution on a seed instruction, yielding five instructions (original + evolved). These instructions are ranked by quality (e.g.,  $C > A = E > B > D$ ) using GPT-4. Based on the rankings, we train the Instruction Ranking Model (IRM) to assess instruction quality. In Figure 1, this process is shown in the left-central segment: “A” represents the original instruction, while “B,” “C,” “D,” and “E” denote the evolved instructions. The first blue arrow illustrates the ranking process via GPT-4, the second arrow shows the ranking outcomes, and the third arrow highlights the use of this ranked data to train the IRM.
- **PRM:** In the middle-right section of Figure 1, the process for training the PRM is depicted. The SFT model generates step-by-step solutions from the given instructions, which are then evaluated and labeled by GPT-4. This labeled data is subsequently used to train the PRM.

### 3. Reinforcement Learning with PPO

In the final step, we integrate the IRM and PRM within a **Proximal Policy Optimization (PPO)**-based reinforcement learning framework. As depicted in the far-right section of Figure 1, the process is as follows:

- (a) The first blue arrow represents instruction scoring by the IRM.
- (b) The second blue arrow shows PPO initialization and the start of reinforcement.
- (c) The third blue arrow illustrates the policy model generating responses based on instructions.
- (d) The fourth blue arrow shows the scoring of each response step using the PRM.
- (e) Arrows five through eight depict the combination of IRM and PRM scores to calculate the final reward score.
- (f) The ninth blue arrow highlights the use of the reward score for the PPO training.

By integrating instruction evolution and reward-based optimization, the RLEIF method significantly enhances the reasoning capabilities of large language models. We hope this explanation resolves some ambiguities and provides a clearer understanding of Figure 1. Thank you again for your valuable suggestions, which will guide us in improving the presentation and clarity of our work.

**Q5: There are also several typos and poorly-worded sentences.**

We sincerely appreciate your thorough review and the time you dedicated to identifying these typos and poorly-worded sentences in our paper. Your attention to detail has been very invaluable. We have corrected all the issues you highlighted in our latest upload of revised paper.

## C.4.3 WEAKNESSES-3

**Somewhat marginal contribution – Evol-Instruct previously existed. PRM previously existed. This paper basically took Evol-Instruct and PRM and used them to train a model. To nitpick a bit, I think a more comprehensive paper would cover more domains such as code.**

**Q1: Somewhat marginal contribution – Evol-Instruct previously existed. PRM previously existed. This paper basically took Evol-Instruct and PRM and used them to train a model**

We sincerely appreciate your valuable feedback. We highlight the key contributions of our paper as follows:

**1. Unlike WizardLM/WizardCoder, which primarily focus on increasing instruction difficulty, we are the first to propose the novel concept of downward evolution, a major distinction in instruction evolution.**

In Table 6 (lines 397–413) of our paper, we provide a detailed analysis of the effects of downward evolution. Specifically, two rounds of downward evolution led to a remarkable improvement in GSM8k performance by 14.8% (74.5 vs. 59.7) and in MATH performance by 19.6% (34.7 vs. 15.1) compared to the original, significantly enhancing the model’s mathematical reasoning capabilities.

Furthermore, our Math Evol-Instruct method outperforms the general evol-instruct approach employed by WizardLM, as elaborated in Appendix A.5 (lines 1608–1629).

This demonstrates that Math Evol-Instruct is instrumental in significantly boosting the model’s mathematical reasoning ability, **as you kindly acknowledged in Strength 1 above.**

**2. In reinforcement learning (RL) training, we firstly propose the instruction quality scoring reward model (IRM) combined with the process supervision reward model (PRM) further enhancing WizardMath mathematical reasoning ability.** As demonstrated in Table 3 (lines 325–336, lines 370-380) of our paper, this approach achieves a remarkable 5%–8% improvement in GSM8k and MATH performance over the SFT backbone across models of various sizes, leveraging PRM and IRM for the PPO training.

**3. We first propose to use AI to annotate the step-level PRM training data.** Additionally, the training datasets for SFT, PRM, and IRM are fully synthesized using AI systems. This fully AI-automated data generation pipeline ensures scalability, **as highlighted in Strength 3 of your feedback.**

**4. WizardMath demonstrates outstanding performance across a wide range of model scales, from 100M to 1B and 70B parameters, on benchmarks such as GSM8k, MATH, and out-of-distribution (OOD) tasks like MWPBench.** It surpasses all existing open-source state-of-the-art models, showcasing the effectiveness and robustness of the RLEIF approach proposed in our study, **as you recognized in Strength 1 above.**

**Q2: I think a more comprehensive paper would cover more domains such as code.**

We sincerely appreciate your insightful suggestions. To explore the effectiveness of our proposed RLEIF method in more other domains such as Code, we replicated the code evol-instruct specifically proposed by WizardCoder for code-related tasks during the SFT stage, and further optimized the PRM step-level label prompts to enhance its compatibility with GPT-4 for annotating PRM training data. Additionally, we compared the performance of ORM and PRM during PPO training. We utilized CodeLlama-Python 7B and 34B as the base models

As shown in the table 34 below, the results demonstrate that for both the CodeLlama-Python 7B and 34B models, Our-Coder-SFT achieved comparable performance to WizardCoder on the HumanEval and MBPP benchmarks. During the PPO training phase, when using CodeLlama-Python 7B as the base model, Our-Coder-RL-PRM showed a 4%-5% improvement on HumanEval and MBPP over Our-Coder-SFT, and significantly outperformed the 2%-3% improvement achieved by Our-Coder-RL-ORM. Similarly, with CodeLlama-Python 34B as the base model, Our-Coder-RL-PRM shows approximately a 4% improvement over Our-Coder-SFT on HumanEval and MBPP, outperforming the 2%-3% improvement of Our-Coder-RL-ORM. These findings underscore the effectiveness of PRM in PPO training for code-related tasks.

In future camera-ready version of our paper, we commit to conducting comprehensive comparisons across more code benchmarks and a broader range of baseline models to further validate the effectiveness of the proposed RLEIF approach.

Table 34: Explore the effectiveness of the proposed RLEIF pipeline in more other domains such as Code

Models	Base	Params	HumanEval	MBPP
CodeLlama-Python-7B				
CodeLlama-Python	-	7B	37.8	57.6
WizardCoder	CodeLlama-Python	7B	48.2	56.6
Our-Coder-SFT	CodeLlama-Python	7B	49.0	56.2
Our-Coder-RL-ORM	CodeLlama-Python	7B	50.5	58.1
Our-Coder-RL-PRM	CodeLlama-Python	7B	53.5	60.4
CodeLlama-Python-34B				
CodeLlama-Python	-	34B	51.8	67.2
WizardCoder	CodeLlama-Python	34B	73.2	73.2
Our-Coder-SFT	CodeLlama-Python	34B	72.7	72.3
Our-Coder-RL-ORM	CodeLlama-Python	34B	74.5	73.7
Our-Coder-RL-PRM	CodeLlama-Python	34B	76.8	76.2

#### C.4.4 QUESTIONS-1

**I find Figure 1 confusing. Why is there a pyramid in the top left and why is it pointing to a pie chart, cube, etc? What are these supposed to be showing? I feel like I am not understanding much from this figure.**

We sincerely apologize for any confusion or inconvenience caused by the current presentation of Figure 1. The pyramid in the top-left corner represents the original seed instructions, while the pie chart, cube, and other icons symbolize the evolved instructions generated through the Math Evol-Instruct method, encompassing both upward and downward evolution. A detailed explanation of the training process depicted in Figure 1 was provided in our response to Weaknesses-2-Q4 above, which we hope will help clarify any uncertainties. In the future camera-ready version of our paper, we will provide a more detailed explanation of Figure 1 to ensure its clarity and comprehensibility. We greatly appreciate your understanding and patience.

#### C.4.5 QUESTIONS-2

I feel like there’s a few missing entries in Table 1. For example, Table 1 shows the results for WizardMath-Mathstral and WizardMath-Qwen2.5, but the base scores of these base models are not shown in the table, so the readers don’t really know how much improvement there is.

We sincerely appreciate your constructive feedback. The table below 35 supplements the performance comparison of Mathstral-7B-v0.1-Base, Qwen2.5-7B-Base, Qwen2.5-Math-1.5B-Base, and Qwen2.5-Math-7B-Base on the GSM8k and MATH datasets.

The results demonstrate that using Mathstral-7B-v0.1-Base as the base model, WizardMath-Mathstral improves performance by 16.7% on GSM8k (93.8 vs. 77.1) and 14.5% on MATH (70.9 vs. 56.6).

When employing Qwen2.5-Math-1.5B-Base as the base model, WizardMath-Qwen2.5-Math-1.5B achieves 9.9% improvement on GSM8k (86.7 vs. 76.8) and 18.8% on MATH (68.6 vs. 49.8).

Similarly, with Qwen2.5-Math-7B-Base, WizardMath-Qwen2.5-Math-7B shows a 2.3% increase on GSM8k (93.9 vs. 91.6) and 22.4% on MATH (77.8 vs. 55.4).

Finally, using Qwen2.5-7B-Base as the base model, WizardMath-Qwen2.5-7B improves by 8.6% on GSM8k (94.0 vs. 85.4) and 24.7% on MATH (74.5 vs. 49.8).



Notably, both Mathstral-7B-v0.1-Base and Qwen2.5-Math-Base, pre-trained on extensive mathematical corpora, exhibit robust mathematical reasoning capabilities and deliver strong performance on GSM8k and MATH datasets. However, our proposed RLEIF method achieves substantial performance enhancements even with these highly math-optimized models. Specifically, on the MATH, RLEIF delivers a performance boost of 15% 25%, while on GSM8k, the improvement ranges from 8% 16% (with the exception of Qwen2.5-Math-7B-Base, which achieves a high baseline of 91.6 on GSM8k but still benefits from a 2.3% enhancement). These results underscore the continuous improvement enabled by our RLEIF method on models pre-trained with specialized mathematical corpora, further validating its effectiveness and scalability.

Table 35: The performance of WizardMath on the GSM8k and MATH based on the Mathstral-7B-v0.1-Base, Qwen2.5-7B-Base, Qwen2.5-Math-1.5B-Base, and Qwen2.5-Math-7B-Base

Models	Base	Params	GSM8k	MATH
Mathstral-v0.1-Base	-	7B	77.1	56.6
WizardMath-Mathstral	Mathstral-v0.1-Base	7B	93.8	70.9
Qwen2.5-Math-Base	-	1.5B	76.8	49.8
WizardMath-Qwen2.5-Math	Qwen2.5-Math-Base	1.5B	86.7	68.6
Qwen2.5-Math-Base	-	7B	91.6	55.4
WizardMath-Qwen2.5-Math	Qwen2.5-Math-Base	7B	93.9	77.8
Qwen2.5-Base	-	7B	85.4	49.8
WizardMath-Qwen2.5	Qwen2.5-Base	7B	94.0	74.5

#### C.4.6 TYPOS:

- L45: struggles -> struggle
- L80: "in recent"
- L91: should say what is IRM first before using the acronym.
- L91: We -> capitalization
- L93: later -> latter
- L105: Jiang et al mentioned without a model (should be mistral?)
- L107: as following -> as follows
- L140,143: spacing
- L145: should be Reinforcement Learning for Large Language Models instead of the other way around?
- L487: "reasoing"

We sincerely appreciate your effort in identifying these typos and poorly-worded sentences in our paper, as well as your thorough and thoughtful review. All of these "Typos" have been carefully addressed and corrected in our latest upload of revised paper.