

WARRIORMATH: ENHANCING THE MATHEMATICAL ABILITY OF LARGE LANGUAGE MODELS WITH A DEFECT-AWARE FRAMEWORK

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) excel in solving mathematical problems, yet their performance is often limited by the availability of high-quality, diverse training data. Existing methods focus on augmenting datasets through rephrasing or difficulty progression but overlook the specific failure modes of LLMs. This results in synthetic questions that the model can already solve, providing minimal performance gains. To address this, we propose WarriorMath, a defect-aware framework for mathematical problem solving that integrates both targeted data synthesis and progressive training. In the synthesis stage, we employ multiple expert LLMs in a collaborative process to generate, critique, and refine problems. Questions that base LLMs fail to solve are identified and iteratively improved through expert-level feedback, producing high-quality, defect-aware training data. In the training stage, we introduce a progressive learning framework that iteratively fine-tunes the model using increasingly challenging data tailored to its weaknesses. Experiments on six mathematical benchmarks show that WarriorMath outperforms strong baselines by 12.57% on average, setting a new state-of-the-art. Our results demonstrate the effectiveness of a defect-aware, multi-expert framework for improving mathematical ability.

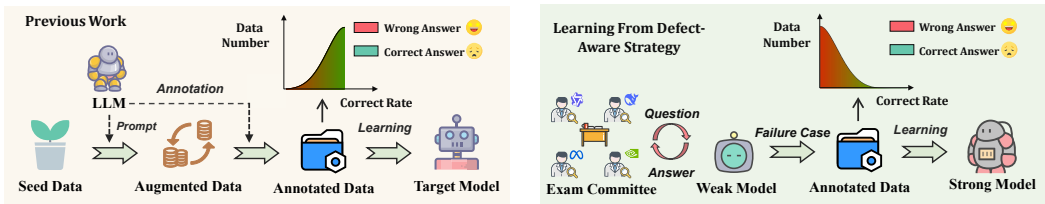
1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable capabilities in solving mathematical and scientific problems (Jaech et al., 2024; Anthropic, 2025; He et al., 2024b; Zeng et al., 2024; OpenAI, 2025; DeepMind, 2025; DeepSeek-AI et al., 2025), positioning them as valuable mathematical assistants. Consequently, enhancing their mathematical ability has become a key research goal. Yang et al. (2024) improve math skills through large-scale pre-training on math data, while Muennighoff et al. (2025); Wen et al. (2025); Min et al. (2024) focus on fine-tuning with high-quality instruction datasets. However, both approaches heavily rely on high-quality data (Xu et al., 2024; Feng et al., 2025), and key training data for strong models (e.g., OpenAI o1 (OpenAI, 2025)) remain private, limiting reproducibility. Thus, collecting and annotating high-quality math problems at scale remains a major bottleneck.

In response, recent research has explored large-scale data synthesis using LLMs to augment training datasets (Tang et al., 2024; Huang et al., 2025; Yue et al., 2024; Liu et al., 2024; Zhou et al., 2024; 2025; Li et al., 2024a; He et al., 2025b; Luo et al., 2025a; Li et al., 2024b; Mei et al., 2025). Some of these methods focus on mining instruction data from pretraining corpora (Yue et al., 2023; Li et al., 2024d), while others generate new data by rephrasing existing problems (Yu et al., 2023) or employing difficulty progression techniques (Xu et al., 2024; Luo et al., 2025a).

While these data synthesis approaches have made progress, they are not tailored to address the specific deficiencies of the base LLM. As a result, they predominantly generate easily solvable problems that do not effectively challenge the model, offering limited learning benefits (see Figure 1). *Inherent defects*, such as misinterpretations of quantifiers or incorrect symbolic steps (Pan et al., 2025; An et al., 2024), are often overlooked. Unlike difficulty, defects reflect internal model limitations. By ignoring these defects, existing synthesis methods primarily generate problems that the model can already solve, hindering further improvement (Yu et al., 2025; Pan et al., 2025; An et al., 2024).

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



(a) Existing data synthesis strategy. Seed datasets are collected, and an external LLM is prompted to augment or label them. However, the resulting problems are often too simple, leading to limited improvements in model performance. (b) Our defect-aware strategy. An Exam Committee of multiple expert LLMs generates, critiques, and refines problems, retaining only the failure data where the base LLMs struggled. This ensures synthesized problems challenge the base LLMs and enhance their capabilities.

Figure 1: Comparisons between our method and traditional data synthesis strategies.

Recent research has shown that aligning training strategies with the model’s evolving defects can significantly improve performance on challenging tasks (An et al., 2024; Wen et al., 2025; Yu et al., 2025; Teams, 2024). Therefore, a more interactive and adaptive synthesis process is required to effectively address these defects.

Based on this insight, we introduce **WarriorMath**, a defect-aware framework for mathematical problem solving that combines both data synthesis and progressive training. As shown in Figure 2, WarriorMath decomposes data synthesis into two stages: **(1) Defect-aware problem construction.** We assemble an exam committee of state-of-the-art mathematical LLMs, each contributing its expertise to generate challenging problems. These problems are then reviewed by judges, while the base LLM provides feedback, identifying the specific areas it has yet to master. **(2) Answer generation and refinement.** Each model attempts to solve the generated problems, and the solutions are filtered and ranked based on Elo ratings and voting. This ensures that the solutions produced are not only correct but also provide meaningful guidance for the base LLM’s learning. Unlike previous methods that expand existing datasets, WarriorMath generates novel, defect-specific training examples from scratch, allowing for more efficient and targeted model improvement.

For training, we introduce a **progressive learning** framework that systematically addresses the model’s defects through a two-stage process. WarriorMath begins with supervised fine-tuning (SFT) using answers generated by multiple expert models, thus establishing a broad foundation of mathematical knowledge. We then identify the failure problems that the model still struggles with and fine-tune it to prioritize stronger solutions on these problematic examples. This process is repeated iteratively, allowing the model to progressively correct its inherent defects without overwriting previously mastered concepts, thereby facilitating more effective learning from its mistakes.

To evaluate the effectiveness of WarriorMath, we conduct evaluations on six prevalent mathematical benchmarks (AIME-2024, 2024; AIME-2025, 2025; AMC-2023, 2023; Lightman et al., 2024; Lewkowycz et al., 2022; He et al., 2024a). Evaluation on these benchmarks indicates that WarriorMath achieves SOTA performance, surpassing existing same-sized open-source large models by an average of 12.57%. Notably, the ablation experiments demonstrate that the proposed synthesis strategy can indeed generate proper high-quality data for the base LLM, as well as the effectiveness of the proposed training framework in learning from defects.

The key contributions of this work include:

- We propose a defect-aware data synthesis pipeline, which generate proper and high-quality data designed for base LLMs from scratch which emulates the educational philosophy of teaching according to aptitude.
- We introduce a progressive learning framework that first learns broadly from experts and then improve ability through iterative alignment focusing on reinforcing knowledge where the model fails while bypassing mastered knowledge.
- We demonstrate that **WarriorMath** achieves state-of-the-art performance among open-source LLMs, with strong data efficiency and generalization, validating the effectiveness of our approach.

2 RELATED WORK

2.1 MATH LLMs

Recent advances in LLMs’ mathematical capabilities have drawn growing attention from both academic and industrial communities. Early successes were fueled by the creation of large-scale pretraining corpora and curated fine-tuning datasets (Paster et al., 2023; Wang et al., 2024; Shao et al., 2024; Yue et al., 2023), which significantly improved model accuracy on standard math benchmarks. This progress was further accelerated by specialized prompting strategies (Wei et al., 2022; Imani et al., 2023), tool augmentation (Gao et al., 2023; Schick et al., 2024), and reinforcement learning techniques (DeepSeek-AI et al., 2025; Zhao et al., 2024). While advanced prompting and test-time scaling methods (Wu et al., 2024; Muennighoff et al., 2025) continue to push performance limits, current LLMs still lag behind those of proprietary ones (e.g. GPT-4o, Claude, etc.), primarily because stronger models often keep their training data proprietary (Hui et al., 2024). As a result, the lack of publicly available high-quality, diverse datasets remains a significant barrier to further development in this field.

2.2 DATA SYNTHESIS

Synthetic data has been employed to augment training datasets for various mathematical LLMs (Luo et al., 2025a; Teams, 2024; Li et al., 2024a). Early approaches follow the Self-Instruct paradigm (Wang et al., 2023), using few-shot prompting to generate synthetic instructions. To enhance diversity and difficulty, recent work (Luo et al., 2025a; An et al., 2024; Liu et al., 2024) further explores iterative refinement and instruction evolution based on reasoning trajectories. While these methods improve the quality and diversity of synthetic data, they primarily focus on difficulty progression rather than addressing the model’s actual capability defects. That is, increasing difficulty does not necessarily target the failure cases of base LLMs. As a result, many synthesized questions remain solvable by the model, providing limited value for improving model ability. Moreover, many of these approaches depend heavily on proprietary LLMs (e.g., GPT-4, Claude) (Muennighoff et al., 2025), making large-scale data generation costly and less reproducible. Recent work such as Feng et al. (2025) proposes a novel paradigm that distills data through multi-agent competitions among open LLMs, reducing reliance on external APIs. However, these methods still overlook the importance of targeting model-specific defects during synthesis.

2.3 LEARNING FROM DEFECT

An emerging line of work investigates how failure cases can be leveraged to guide LLMs toward improved performance. Reflexion (Shinn et al., 2023) introduces a self-reflection mechanism in which the model analyzes past failures using either internal reflections or external feedback. Similarly, Gou et al. (2024) use external tools to provide real-time critiques, while Chen et al. (2024) enable models to execute and debug code to enhance factual consistency. In contrast to these tool-based feedback approaches, WizardLM-2 (Teams, 2024) proposes the AI-Align-AI (AAA) framework, where multiple LLMs collaboratively teach and critique each other. This setup involves simulated dialogues, quality evaluations, and constructive suggestions for improvement, allowing models to iteratively refine their outputs in a multi-turn process. Despite these advances, most methods still treat failure feedback as a post-hoc augmentation rather than an integral part of the training process. Furthermore, few approaches systematically incorporate failure-driven supervision into dataset synthesis pipelines, leaving a gap between training data construction and the model’s evolving weaknesses.

3 METHOD

In this section, we elaborate on the details of our WarriorMath. As illustrated in Figure 2, the pipeline mainly contains two components: Defect-aware data synthesis and Progressive training. The details of data synthesis will be presented in § 3.1 and the training method will be described in § 3.2

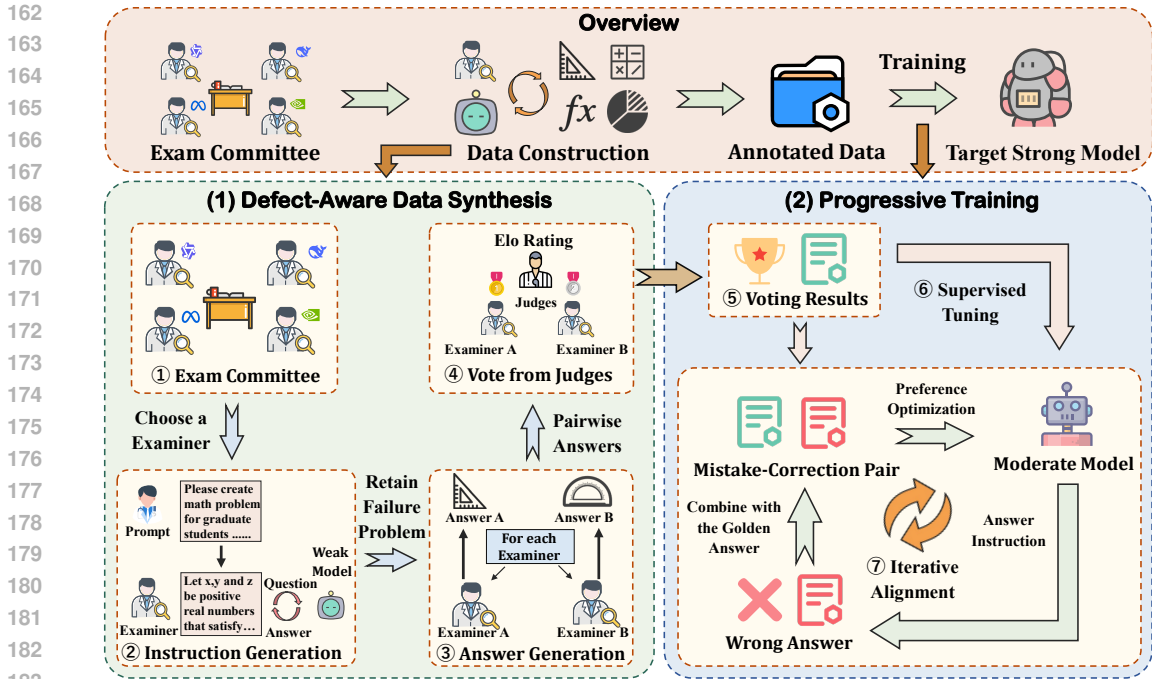


Figure 2: **The WarriorMath’s pipeline.** The WarriorMath pipeline includes two key components: **Defect-Aware data synthesis** and **Progressive Training**. (1) A multi-expert committee generates problems via the Examiner. The Base Model identifies its own weaknesses by attempting these problems, retaining only the failed problems. Solvers generate competing solutions, and Judges use Elo ratings to select the authoritative Answer, constructing the training data. (2) The model undergoes Supervised Fine-Tuning (SFT). It then enters the Iterative Defect Alignment loop, where the model’s failed responses (hard negatives) are paired with the Golden Answers and used for Preference Optimization (DPO-based) to systematically overcome the identified defects.

3.1 DEFECT-AWARE DATA SYNTHESIS

Unlike conventional methods that expand existing datasets, WarriorMath generates novel, defect-specific training examples from scratch. We employ a collaborative **Exam Committee** of expert LLMs to generate, solve, and critique mathematical problems. To ensure clarity in our multi-agent framework, we explicitly define three distinct roles for the committee members: the **Examiner** (generates the problem), the **Solver** (proposes solutions), and the **Judge** (evaluates solution quality).

3.1.1 PROBLEM GENERATION

Committee members Setting. The effectiveness of **WarriorMath** is directly linked to the strength and diversity of its members. For this study, we select five leading open-source math LLMs, DeepSeek-R1-Distill-Llama-70B (DeepSeek-AI et al., 2025), Qwen2.5-Math-72B-Instruct (Yang et al., 2024), QwQ-32B (Qwen, 2025), AceMath-72B-Instruct (Liu et al., 2025b) and Phi-4-reasoning (Abdin et al., 2025). In each synthesis round, one model is selected as the Examiner, while the others act as Solvers or Judges.

Problem Synthesis from Scratch The goal is to leverage the Examiner’s capabilities to pose challenging questions to expose the defects of the **Base Model**. Unlike previous approaches where problems are generated indiscriminately, our design ensures the Examiner focuses on generating problems that the current Base Model is likely to fail. We design a prompt that instructs the Examiner to act as a world-class expert generating difficult mathematical problems. To achieve this, we design a prompt that instructs LLM A to act as a world-class expert in generating difficult and diverse mathematical problems. Specifically, the prompt emphasizes four dimensions: (1) Quality, requiring

problems to be clear, well-structured, and unambiguous; (2) Difficulty, demanding deep mathematical reasoning beyond pattern matching; (3) Diversity, covering a broad set of domains such as algebra, calculus, discrete mathematics, and geometry; and (4) Challenge, encouraging adversarial elements like subtle traps or misleading intermediate steps. This deliberate prompt design ensures LLM A activates its mathematical capabilities while producing maximally diagnostic questions to evaluate LLM B. Further Problem Synthesis details are provided in Appendix B.

Deduplication and Defect-Aware assessment. To ensure the quality and informativeness of selected problems, we identify and filter out problems that are repetitive, ambiguous, or excessively difficult for the base LLM to learn from effectively. The Prompt for selection and discarded problems are provided in appendix B.3. To further ensure that the retained problems are truly valuable for learning, we conduct a defect-aware assessment of the model’s performance. For each problem, we generate $N = 16$ rollouts using the base model. After verifying the correctness of these outputs, we discard problems for which all sampled solutions are correct, retaining only those that the model answers incorrectly. Through interactive assessment, we confirm that these retained problems continue to provide meaningful learning signals. Finally, to preserve both the diversity and representativeness of the instruction set, we apply the KCenterGreedy algorithm (Sener & Savarese, 2018) to select a final subset \bar{I} , using the *all-roberta-large-v1* embedding model (Liu et al., 2019) to compute semantic similarity between instructions.

3.1.2 ANSWER GENERATION AND REFINEMENT

Once a defect-aware problem is identified, we should generate a high-quality “Golden Answer.” This is achieved through a competitive refinement phase involving Solvers and Judges.

Answer Generation (The Solvers). For every instruction i , we assign two models to the role of Solver. Solver A (The Examiner): The model that created the problem must provide its own solution. Solver B (Peer Expert): A different committee member provides an independent solution. This generates a pairwise set of answers for evaluation. **Answer Refinement (The Judges).** The remaining committee members act as Judges, voting on the correctness and helpfulness of the paired responses (more details about can be found in AppendixB.4). We first calculate the *local score* based on the raw vote counts (t_A, t_B)

$$x_{A>B}^i = \frac{t_A}{t_A + t_B} \quad x_{B>A}^i = \frac{t_B}{t_A + t_B} \quad (1)$$

Here, $x_{A>B}^i$ represents the percentage of votes Solver A receives, while $x_{B>A}^i$ similarly represents the percentage of votes Solver B receives. t_A and t_B are the raw vote counts for A and B.

Elo Rating Integration. Relying solely on local votes can be problematic due to judge bias or stochasticity, potentially allowing weaker models to win unrepresentatively. To address this limitation and enforce **global consistency**, We introduce the concept of *the Elo rating* Bai et al. (2022), which provides a more comprehensive reflection of a model’s relative performance over time and across various evaluations.

$$X_{A>B}^{Elo} = \frac{1}{1 + 10^{(R_B - R_A)/400}} \quad (2)$$

$$X_{B>A}^{Elo} = \frac{1}{1 + 10^{(R_A - R_B)/400}}$$

where $X_{A>B}^{Elo}$ and $X_{B>A}^{Elo}$ indicate the expected probabilities of A defeating B and B defeating A, respectively. R_A and R_B are the Elo rating of A and B, which are updated dynamically and iteratively. The update rule is:

$$R_A \leftarrow R_A + K \times (s_{A>B}^i - X_{A>B}^{Elo}) \quad (3)$$

$$R_B \leftarrow R_B + K \times (s_{B>A}^i - X_{B>A}^{Elo})$$

where s^i is the actual outcome (1 for win, 0.5 for draw, 0 for loss) and K controls sensitivity.

Based on Equation 1 and Equation 2, we calculate the definitive score for Solver A’s response to instruction i by balancing the local vote with the global Elo expectation:

$$e_A^i = \sum_{B \in Com \setminus A} \alpha X_{A>B}^{Elo} + (1 - \alpha)x_{A>B}^i \quad (4)$$

where Com is the set of all the solvers and ‘\’ is the subtraction operation. α is the coefficient to balance the local contingency and global consistency.

3.2 PROGRESSIVE TRAINING

We introduce a **Progressive Training** framework that incrementally improves a model’s mathematical reasoning by systematically identifying and correcting its errors. The framework consists of two stages: supervised fine-tuning (SFT) and iterative alignment.

Stage 1: Supervised Fine-Tuning. Given a dataset $\mathcal{D} = \{(x_i, Y_i, \{r_i^j\}_{j=1}^N)\}_{i=1}^N$, where x_i is an instruction, $Y_i = \{y_i^j\}_{j=1}^N$ are expert responses, and r_i^j are their scores, we select the highest-scoring response as $y_i^{\text{gold}} = \arg \max_{y_i^j \in Y_i} r_i^j$. The gold pairs (x_i, y_i^{gold}) are then used to initialize the model M_0 through maximum likelihood estimation.

Stage 2: Iterative Alignment. To further refine the model, we adopt the iterative alignment strategy of Pang et al. (2024), where the model improves by learning from its own mistakes. At iteration t , the current model M_t generates a set of candidate responses

$$G_i = \{(c_i^n, y_i^n)\}_{n=1}^{N_i} \sim M_t(x_i), \quad (5)$$

with reasoning traces c_i^n and final answers y_i^n . Each response is labeled for correctness by $r_i^n = \mathcal{R}(y_i^n, \hat{y}_i)$, which reduces to $r_i^n = 1$ if $y_i^n = \hat{y}_i$, and 0 otherwise. This yields a labeled set $G_i = \{(c_i^n, y_i^n, r_i^n)\}$ and the subset of incorrect responses

$$G_i^{\text{neg}} = \{(c_i^n, y_i^n) \mid r_i^n = 0\}. \quad (6)$$

For each gold pair $(c_i^{\text{gold}}, y_i^{\text{gold}})$, we construct preference pairs by contrasting it with negatives $(c_i^l, y_i^l) \in G_i^{\text{neg}}$, forming D_t^{pairs} . The model is optimized with a hybrid objective:

$$\begin{aligned} \mathcal{L}_{\text{total}} &= \mathcal{L}_{\text{DPO}} + \alpha \mathcal{L}_{\text{NLL}}, \\ \mathcal{L}_{\text{DPO}} &= -\log \sigma \left(\beta \log \frac{M_\theta(c_i^{\text{gold}}, y_i^{\text{gold}} | x_i)}{M_t(c_i^{\text{gold}}, y_i^{\text{gold}} | x_i)} \right. \\ &\quad \left. - \beta \log \frac{M_\theta(c_i^l, y_i^l | x_i)}{M_t(c_i^l, y_i^l | x_i)} \right), \\ \mathcal{L}_{\text{NLL}} &= -\frac{1}{|c_i^{\text{gold}}| + |y_i^{\text{gold}}|} \log M_\theta(c_i^{\text{gold}}, y_i^{\text{gold}} | x_i). \end{aligned} \quad (7)$$

Here σ is the sigmoid, and α, β are hyperparameters. After each optimization step, the updated model $M_{t+1} = M_\theta$ is used to generate new responses for the next iteration.

Through repeated self-correction, the model is progressively aligned with expert reasoning, thereby accumulating the collective expertise of diverse mathematical committee members.

4 EXPERIMENT

4.1 EXPERIMENTAL SETUP

Backbones. We implement WarriorMath with two initialization backbones: (1) WarriorMath-Qwen, initialized from Qwen2.5-Math-7B Yang et al. (2024); (2) WarriorMath-DS, initialized from DeepSeek-R1-Distill-Qwen-7B (DeepSeek-AI et al., 2025). As for the competitors of expert battles, we choose strong open-source LLMs including DeepSeek-R1-Distill-Llama-70B (DeepSeek-AI et al., 2025), Qwen2.5-Math-72B-Instruct (Yang et al., 2024), QwQ-32B (Qwen, 2025), AceMath-72B-Instruct (Liu et al., 2025b), Phi-4-reasoning (Abdin et al., 2025).

Table 1: Evaluation results across six mathematical reasoning benchmarks. We report Pass@1 accuracy (mean \pm std) of all models across six math benchmarks under a standardized evaluation setup—results are averaged over ten seeds for AIME and AMC, and three seeds for the rest.

Models	Base	AIME'24	AIME'25	AMC'23	MATH500	Minerva	Olympiad
<i>Expert Teacher Models</i>							
Qwen2.5-Math-72B-Instruct (Qwen et al., 2025)	Qwen2.5-Math-72B	32.0 \pm 5.9	26.3 \pm 7.3	59.7 \pm 6.1	85.2 \pm 0.5	44.1 \pm 2.2	49.0 \pm 0.5
AceMath-72B-Instruct (Liu et al., 2025b)	Qwen2.5-Math-72B	31.3 \pm 7.7	28.8 \pm 2.2	60.1 \pm 2.4	86.1 \pm 2.3	57.0 \pm 1.6	48.4 \pm 1.3
DeepSeek-R1-Distill-Llama-70B (DeepSeek-AI et al., 2025)	Llama-3-70B	67.0 \pm 1.9	55.3 \pm 5.7	96.8 \pm 2.1	95.1 \pm 0.7	45.1 \pm 1.7	73.8 \pm 0.5
QwQ-32B (Qwen, 2025)	Qwen2.5-32B	76.3 \pm 3.3	69.0 \pm 4.5	96.2 \pm 2.4	97.5 \pm 0.6	49.0 \pm 0.2	78.1 \pm 1.0
Phi-4-reasoning (Abdin et al., 2025)	Phi-4-14B	74.6 \pm 5.1	63.1 \pm 6.3	96.0 \pm 2.7	97.0 \pm 0.3	49.8 \pm 0.2	78.0 \pm 0.8
<i>Qwen-Based Models</i>							
Qwen2.5-Math-7B-Base Yang et al. (2024)	-	20.7 \pm 3.8	8.7 \pm 3.9	56.2 \pm 5.7	64.3 \pm 0.5	17.3 \pm 1.9	29.0 \pm 0.5
Qwen2.5-Math-7B-Instruct (Yang et al., 2024)	Qwen2.5-Math-7B	15.7 \pm 3.9	10.7 \pm 3.8	67.0 \pm 3.9	82.9 \pm 0.1	35.0 \pm 0.6	41.3 \pm 0.9
Qwen-2.5-Math-7B-SimpleRL-Zoo (Zeng et al., 2025)	Qwen2.5-Math-7B	22.7 \pm 5.2	10.7 \pm 3.4	62.2 \pm 3.6	76.9 \pm 1.8	30.1 \pm 2.8	39.3 \pm 0.6
Qwen2.5-Math-7B-Oat-Zero (Liu et al., 2025a)	Qwen2.5-Math-7B	28.0 \pm 3.1	8.8 \pm 2.5	66.2 \pm 3.6	79.4 \pm 0.3	34.4 \pm 1.4	43.8 \pm 1.1
LIMR (Li et al., 2025)	Qwen2.5-Math-7B	30.7 \pm 3.2	7.8 \pm 3.3	62.2 \pm 3.4	76.5 \pm 0.4	34.9 \pm 1.3	39.3 \pm 0.9
Qwen2.5-7B-Instruct (Qwen et al., 2025)	Qwen2.5-7B	12.3 \pm 3.2	7.3 \pm 3.4	52.8 \pm 4.8	77.1 \pm 1.2	34.9 \pm 1.0	38.7 \pm 1.4
s1.1-7B (Muennighoff et al., 2025)	Qwen2.5-7B	19.0 \pm 3.2	21.0 \pm 5.5	59.5 \pm 3.7	80.8 \pm 0.6	37.5 \pm 1.1	48.2 \pm 1.4
Eurus-2-7B-PRIME (Cui et al., 2025)	Qwen2.5-7B	17.8 \pm 2.2	14.0 \pm 1.7	63.0 \pm 3.9	80.1 \pm 0.1	37.5 \pm 1.0	43.9 \pm 0.3
Bespoke-Stratos-7B (Bespoke Labs, 2024)	Qwen2.5-7B	20.3 \pm 4.3	18.0 \pm 4.8	60.2 \pm 4.9	84.7 \pm 0.5	39.1 \pm 1.3	51.9 \pm 1.1
WarriorMath-Qwen-7B	Qwen2.5-7B	48.3\pm2.6	36.5\pm5.0	83.0\pm2.5	88.3\pm1.4	41.2\pm2.8	52.1\pm0.8
<i>DeepSeek-Based Models</i>							
DeepSeek-R1-Distill-Qwen-7B (DeepSeek-AI et al., 2025)	-	52.3 \pm 6.3	39.0 \pm 5.9	91.5 \pm 2.7	94.1 \pm 0.3	40.1 \pm 0.4	67.3 \pm 0.1
Light-R1-DS-7B (Wen et al., 2025)	DeepSeek-R1-Distill-Qwen-7B	53.0 \pm 4.8	41.0 \pm 3.5	90.0 \pm 3.1	93.5 \pm 0.5	41.3 \pm 1.3	68.0 \pm 1.2
AReal-boba-RL-7B (inclusionAI, 2025)	DeepSeek-R1-Distill-Qwen-7B	56.7 \pm 9.2	40.0 \pm 9.1	90.0 \pm 4.8	94.4 \pm 1.0	40.8 \pm 3.0	68.4 \pm 1.8
WarriorMath-DS-7B	DeepSeek-R1-Distill-Qwen-7B	60.0\pm9.1	50.7\pm9.1	93.2\pm4.9	95.0\pm1.0	43.20\pm1.9	69.6\pm1.8

Datasets. To evaluate **WarriorMath**'s mathematical capabilities, we use six prevalent benchmarks: (1) **AIME 2024**, **AIME 2025** (AIME-2024, 2024; AIME-2025, 2025) are benchmarks that include particularly challenging math problems from the American Invitational Mathematics Examination (AIME) of 2024 and 2025, designed to assess advanced problem-solving skills. (2) **AMC 2023**: (AMC-2023, 2023) High school-level problems from the American Mathematics Competitions (AMC), testing core mathematical concepts (3) **MATH-500** (Lightman et al., 2024) is a dataset containing high school-level math problems. It serves to assess a model's ability to handle more advanced mathematical reasoning; and (4) **Minerva** (Lewkowycz et al., 2022) is a benchmark dataset comprising a diverse collection of quantitative reasoning problems that cover topics such as arithmetic, algebra, geometry, calculus, physics, and chemistry, with difficulty levels ranging from grade school to college. (5) **Olympiad Bench**: He et al. (2024a) an Olympiad-level bilingual scientific benchmark, featuring 8,476 problems from Olympiad-level mathematics and physics competitions. To ensure accurate evaluation, we follow the evaluation method proposed by Hochlehnert et al. (2025). We use the Pass@1 metric as the primary evaluation criterion. For each result, we report the mean and standard deviation computed over multiple random seeds. All experiments are conducted using `lighteval` (Fourrier et al., 2023) with the `vllm` backend (Kwon et al., 2023).

Baseline. We compare our method against reinforcement learning (RL) approaches trained on the Qwen2.5 Math Base models, including Oat-Zero (Liu et al., 2025a), LIMR (Li et al., 2025), and SimpleRL-Zoo (Zeng et al., 2025). We also evaluate our approach against supervised fine-tuning (SFT) baselines, such as s1.1 (Muennighoff et al., 2025), Eurus2 Prime (Cui et al., 2025), Bespoke Stratos (Bespoke Labs, 2024), OpenR1 (Face, 2025) and OpenThinker (Team, 2025). In addition, we consider recent state-of-the-art methods based on `deepseek-r1-distill-qwen-7b` as the backbone, such as `LightR1` (Wen et al., 2025), and `AReal-boba-RL-7b` (inclusionAI, 2025).

Implementation Details During the data synthesis, we adopt 9 different generation configs where temperature $t \in \{0.60, 0.65, 0.70\}$ and top-p $p \in \{0.85, 0.90, 0.95\}$. The detailed prompts can be found in Appendix B.1. As for the training stage, the global batch size is set to 512, and the number of total training steps is set to 448. We use a learning rate of 1×10^{-5} and a weight decay of 3×10^{-7} . Additionally, a WarmupLR scheduler with a warmup ratio of 0.2 is used.

4.2 PERFORMANCE AND COMPARISON

Main Result. The results on the math benchmarks are summarized in Table 1. **WarriorMath** achieves SOTA performance, with a pass@1 accuracy of 60% in AIME'24 and 56.7% in AIME'25, surpassing all other fine-tuned models. This highlights the efficacy of our approach in generating high-quality data and effective training process.

Table 2: Performance of different data synthesis strategies on three mathematical benchmarks.

Models	Base	GSM8K	MATH-500	AIME2024
Qwen2.5-Math-7B-Instruct Yang et al. (2024)	-	95.2	83.6	13.3
Openmathinstruct-7B Toshniwal et al. (2025)	Qwen2.5-Math-7B	92.0	79.6	10.0
NuminaMath-7B Li et al. (2024c)	Qwen2.5-Math-7B	92.9	81.8	20.0
Evol-Instruct-7B Luo et al. (2025a)	Qwen2.5-Math-7B	88.5	77.4	16.7
KPDDS-7B Huang et al. (2025)	Qwen2.5-Math-7B	89.9	76.0	10.0
PROMPTCOT-Qwen-7B Zhao et al. (2025)	Qwen2.5-Math-7B	93.3	84.0	26.7
WarriorMath-Qwen-7b-SFT	Qwen2.5-7B	95.7	83.8	36.7

Table 3: The proportion of different tasks in the training data.

Mathematics Domain	Percentage (%)	Definition
Applied Mathematics	11.4	Apply mathematical methods to solve cross-field problems.
Algebra	30.3	Study of mathematical symbols and manipulation rules.
Discrete Mathematics	12.9	Study of discrete mathematical structures (e.g., graphs, integers).
Geometry	14.7	Study of properties/relations of points, lines, surfaces, solids.
Number Theory	13.1	Study of integers and integer-valued functions.
Precalculus	1.2	Math preparation for calculus (functions, trigonometry).
Calculus	1.8	Study of continuous change (derivatives, integrals).
Differential Equations	0.5	Equations involving derivatives for quantity change.

Table 4: Results with Iterative Alignment.

Model	AIME’24
WarriorMath-Qwen-7b-SFT	36.7±4.5
Iteration 1	42.5±5.5
Iteration 2	44.7±7.2
Iteration 3	48.3±2.6

Table 5: Results when learning from varying numbers of committee members.

#Num	AIME’24	AIME’25	AMC’23
1	19.7±2.9	15.7±2.7	59.5±4.5
2	29.4±3.2	23.5±4.2	69.3±3.6
5	48.3±2.6	36.5±5.0	83.0±2.5

Data Quality. In order to assess the effectiveness of the problem generation pipeline, we compare our method with the following problem generation baselines: (1) **Evol-Instruct**: This method (Luo et al., 2025a) aims to enhance the quality of instruction data by improving both its complexity and diversity, thus facilitating the generation of more varied and challenging problems; (2) **KPDDS**: A data synthesis framework (Huang et al., 2025) that generates question-answer pairs by leveraging key concepts and exemplar practices derived from authentic data sources; (3) **OpenMathInstruct**: This method (Toshniwal et al., 2025) utilizes few-shot learning to prompt an LLM to create new math problems based on existing examples, without explicit instructions for adjusting difficulty or introducing new constraints; (4) **NuminaMath**: This approach (Li et al., 2024c) uses an LLM to generate novel math questions starting from a reference problem; (5) **PROMPTCOT**: This method (Zhao et al., 2025) synthesizes complex problems based on mathematical concepts and the rationale behind problem construction, emulating the thought processes of experienced problem designers. For fair comparisons we follow the evaluation scripts provided in (Zhao et al., 2025). The results of data quality assessment, presented in Tables 2, Our method achieves state-of-the-art performance across multiple benchmarks, outperforming the baselines, which highlights the efficacy of our defect-aware approach in generating high-quality problems

Iterative Defect Alignment. As a seed model M_0 we use the WarriorMath-Qwen-7b-SFT, which is fine-tuned with instruction data generated by Defect-Aware Committee Assessment. In each iteration, we generate $N = 32$ solutions per problem using sampling with temperature 0.7 and top-p 0.8, and verify the answer to select wrong solution (c_i, y_i) in the loser set G_i^l . Then we generate $K = 10$ pairs per problem for training with our loss in Equation 3.2. In total, we perform three iterations, producing models M_1, M_2, M_3 . The coefficient α is tuned in $\{0.25, 0.5, 1\}$ when training M_1 , and we end up using 1 for all experiments in the paper. The coefficient β in the DPO loss is tuned in $\{0.05, 0.1, 0.5\}$, and we end up using 0.1 in this experiment.

Overall results are given in Table 4. We find that WarriorMath outperforms supervised fine-tuning (SFT) on the gold (dataset-provided) data, and steady growth over the iteration rounds.

Table 6: Ablation results across six mathematical reasoning benchmarks. We report Pass@1 accuracy (mean \pm std) across six reasoning benchmarks, consistent with the setup in Table 1.

Models	AIME'24	AIME'25	AMC'23	MATH500	Minerva	Olympiad
WarriorMath-Qwen-7B-SFT	36.7 \pm 4.5	36.5 \pm 5.0	83.0 \pm 2.5	88.3 \pm 1.4	41.2 \pm 2.8	52.1 \pm 0.8
WarriorMath-Qwen-7B-SFT (w/o defect)	32.0 \pm 5.9	21.4 \pm 4.3	74.1 \pm 5.2	86.6 \pm 0.6	31.7 \pm 0.6	51.9 \pm 0.4
WarriorMath-Qwen-7B-SFT (w/o elo rating)	35.4 \pm 5.6	24.0 \pm 4.1	79.5 \pm 5.1	86.7 \pm 1.6	38.4 \pm 2.2	50.1 \pm 2.3

4.3 ABLATION STUDY.

4.3.1 NUMBER OF COMMITTEE MEMBERS

Table 5 presents the results observed when the target model learns from varying numbers of committee members. The target model shows a significant improvement when learning from just one mathematical LLM, indicating that even a single committee member enables it to acquire a specific set of knowledge. However, as the number of members increases, **WarriorMath** benefits from learning across all mathematical LLMs. As a result, the model trained with 5 math LLMs outperforms others across all 6 benchmarks, demonstrating the advantages of integrating knowledge from multiple specialized experts.

4.3.2 DEFECT FILTERING

Our *defect-aware filtering* specifically targets problems proposed by teacher models, ensuring that the resulting data remain sufficiently challenging for downstream learners (e.g., Qwen2.5-Math, Qwen-DeepSeek-R1-Distill-7B). As base models grow more capable, many mined problems may be trivial; thus, defect-aware filtering becomes indispensable. Recent works (Muennighoff et al., 2025; He et al., 2025a; Zhang et al., 2025) independently adopt performance-based difficulty screening, further validating its effectiveness. To empirically assess the role of *base model*, we ablated the Defect-Aware Committee Assessment stage by removing them. Results in Table 6 clearly demonstrate that Defect Filtering with base model are crucial for optimal performance across benchmarks.

4.3.3 ELO RATING

Table 6 demonstrates the impact of the Judge system. Replacing the judge’s Elo rating with mediate voting drops AIME’25 performance from 36.5% to 21.4%. Mediate voting alone suffer from high variance (noise from single judges), whereas Elo provides global consistency, representing more robust and accurate measure of a model’s overall ability.

4.3.4 DATA ANALYSIS

Data Dependence As discussed in Section 3.1.1, our data sources are derived from multi-expert LLMs, which distinguishes them from the commonly used datasets in existing research. To assess the novelty and dependence of our data, we conducted an analysis across all datasets. Figure 3 illustrates the overlap between the instructions mined from expert LLMs and those in two widely adopted math training datasets: (1) DeepScaleR (Luo et al., 2025b) and (2) Omni-MATH (Gao et al., 2025), measured by the ROUGE score. The majority of the mined instructions show a ROUGE score below 0.3, indicating significant distinctiveness from the existing datasets. Notably, no mined instruction exceeds a ROUGE score of 0.6, further confirming that these instructions are generated from the internal distribution of expert LLMs, rather than being simple replications or extensions of the training data. This unique source ensures a higher degree of independence, making the instructions especially valuable for model training by providing novel examples that can enhance the model’s capabilities.

Data Diversity A key characteristic of our training data is its extensive topical diversity across mathematical domains. As shown in Table 3, the classification results highlight that WarriorMath covers a wide range of core mathematical fields. These include foundational subjects like Applied Mathematics and Basic Geometry, as well as more advanced areas such as Number Theory and Differential Equations. This broad topical range ensures that models trained on WarriorMath are exposed to a rich variety of mathematical concepts and problem-solving strategies, enabling the development of more robust and versatile mathematical capabilities.

5 CONCLUSION

In this work, we present WarriorMath, a defect-aware framework that enhances mathematical ability in LLMs through defect-aware data synthesis and progressive training. Our method constructs high-quality, defect-aware training data by leveraging a committee of expert LLMs to generate, critique, and refine problems specifically designed to expose the base LLM’s inherent defects. Through a two-stage progressive training process, WarriorMath incrementally aligns the model to stronger mathematical ability. Extensive experiments on six mathematical benchmarks demonstrate that WarriorMath achieves state-of-the-art performance among open-source models, highlighting the importance of learning from model-specific defects.

6 LIMITATION

This paper allows for the low-cost generation of high-quality and diverse data from scratch. However, as the number of expert models increases, the evaluation process becomes increasingly time-consuming. Designing more efficient and scalable multi-agent collaboration mechanisms remains an important direction for future research. **Another limitation concerns the coarseness of the defect signal. Our current definition of a defect is a pragmatic, binary signal based solely on the final answer’s correctness. This prevents the framework from providing the most fine-grained feedback.**

ETHICS STATEMENT

This work does not involve human subjects, personally identifiable information, or sensitive data. The datasets used are publicly available, and all experiments comply with the ICLR Code of Ethics.

REPRODUCIBILITY STATEMENT

We have made extensive efforts to ensure reproducibility. The detailed methodology of our proposed approach is presented in Section 3, while the experimental settings, including training procedures and evaluation protocols, are described in Section 4. To further support reproducibility, we plan to release the complete source code and instructions upon the acceptance of this paper.

REFERENCES

- Marah Abdin, Sahaj Agarwal, Ahmed Awadallah, Vidhisha Balachandran, Harkirat Behl, Lingjiao Chen, Gustavo de Rosa, Suriya Gunasekar, Mojan Javaheripi, Neel Joshi, Piero Kauffmann, Yash Lara, Caio César Teodoro Mendes, Arindam Mitra, Besmira Nushi, Dimitris Papailiopoulos, Olli Saarikivi, Shital Shah, Vaishnavi Shrivastava, Vibhav Vineet, Yue Wu, Safoora Yousefi, and Guoqing Zheng. Phi-4-reasoning technical report, 2025. URL <https://arxiv.org/abs/2504.21318>.
- AIME-2024. Aime24, 2024. URL <https://huggingface.co/datasets/math-ai/aime24>.
- AIME-2025. Aime25, 2025. URL <https://huggingface.co/datasets/math-ai/aime25>.
- AMC-2023. Amc23, 2023. URL <https://huggingface.co/datasets/math-ai/amc23>.
- Chenyang An, Zhibo Chen, Qihao Ye, Emily First, Letian Peng, Jiayun Zhang, Zihan Wang, Sorin Lerner, and Jingbo Shang. Learn from failure: Fine-tuning llms with trial-and-error data for intuitionistic propositional logic proving. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of ACL 2024*, pp. 776–790, 2024. URL <https://aclanthology.org/2024.acl-long.45/>.
- Anthropic. Claude 3.7 Sonnet System Card, 2025. URL <https://assets.anthropic.com/m/785e231869ea8b3b/original/claude-3-7-sonnet-system-card.pdf>. Accessed: 2025-03-29.

- 540 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
541 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson
542 Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez,
543 Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario
544 Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan.
545 Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022.
546 URL <https://arxiv.org/abs/2204.05862>.
- 547 Bespoke Labs. Bespoke-stratos-7b. [https://huggingface.co/bespokelabs/
548 Bespoke-Stratos-7B](https://huggingface.co/bespokelabs/Bespoke-Stratos-7B), 2024. Accessed: 2025-03-29.
- 549
550 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models
551 to self-debug. In *Proceedings of ICLR 2024*, 2024. URL [https://openreview.net/forum?id=
552 KuPixIqPiq](https://openreview.net/forum?id=KuPixIqPiq).
- 553
554 Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu
555 Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint
556 arXiv:2502.01456*, 2025. URL <https://arxiv.org/abs/2502.01456>.
- 557
558 Google DeepMind. Gemini 2.5: Our most intelligent ai model, 2025. URL [https://blog.google/
559 technology/google-deepmind/gemini-model-thinking-updates-march-2025/](https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/). Accessed:
560 2025-04-07.
- 561
562 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,
563 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,
564 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao
565 Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
566 Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,
567 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding,
568 Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang
569 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong,
570 Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao,
571 Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang,
572 Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang,
573 Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L.
574 Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang,
575 Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng
576 Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng
577 Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan
578 Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang,
579 Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen,
580 Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li,
581 Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang,
582 Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan,
583 Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia
584 He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong
585 Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha,
586 Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang,
587 Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li,
588 Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen
589 Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
590 URL <https://arxiv.org/abs/2501.12948>.
- 591
592 Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL [https:
593 //github.com/huggingface/open-r1](https://github.com/huggingface/open-r1).
- 594
595 Huawei Feng, Pu Zhao, Qingfeng Sun, Can Xu, Fangkai Yang, Lu Wang, Qianli Ma, Qingwei Lin,
596 Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. Warriorcoder: Learning from expert battles to
597 augment code large language models, 2025. URL <https://arxiv.org/abs/2412.17395>.

- 594 Clémentine Fourier, Nathan Habib, Hynek Kydlíček, Thomas Wolf, and Lewis Tunstall. LightEval:
595 A lightweight framework for LLM evaluation, 2023. URL [https://github.com/huggingface/
596 lighteval](https://github.com/huggingface/lighteval).
- 597
- 598 Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma,
599 Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan,
600 Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. Omni-math:
601 A universal olympiad level mathematic benchmark for large language models. In *Proceedings of
602 ICLR 2025*, 2025. URL <https://openreview.net/forum?id=yaqPf0KA1N>.
- 603 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
604 Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine
605 Learning*, pp. 10764–10799. PMLR, 2023. URL <https://arxiv.org/abs/2211.10435>.
- 606
- 607 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen.
608 Critic: Large language models can self-correct with tool-interactive critiquing, 2024. URL
609 <https://arxiv.org/abs/2305.11738>.
- 610
- 611 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han,
612 Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. OlympiadBench: A
613 challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific
614 problems. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of ACL 2024*,
615 pp. 3828–3850, August 2024a. URL <https://aclanthology.org/2024.acl-long.211/>.
- 616
- 617 Jujie He, Jiakai Liu, Chris Yuhao Liu, Rui Yan, Chaojie Wang, Peng Cheng, Xiaoyu Zhang, Fuxiang
618 Zhang, Jiacheng Xu, Wei Shen, Siyuan Li, Liang Zeng, Tianwen Wei, Cheng Cheng, Bo An,
619 Yang Liu, and Yahui Zhou. Skywork open reasoner 1 technical report, 2025a. URL <https://arxiv.org/abs/2505.22312>.
- 620
- 621 Minghua He, Tong Jia, Chiming Duan, Huaqian Cai, Ying Li, and Gang Huang. Llmelog: An
622 approach for anomaly detection based on llm-enriched log events. In *2024 IEEE 35th International
623 Symposium on Software Reliability Engineering (ISSRE)*, pp. 132–143. IEEE, 2024b.
- 624
- 625 Minghua He, Fangkai Yang, Pu Zhao, Wenjie Yin, Yu Kang, Qingwei Lin, Saravan Rajmohan,
626 Dongmei Zhang, and Qi Zhang. Execoder: Empowering large language models with executability
627 representation for code translation. *arXiv preprint arXiv:2501.18460*, 2025b.
- 628
- 629 Andreas Hochlehnert, Hardik Bhatnagar, Vishaal Udandarao, Samuel Albanie, Ameya Prabhu, and
630 Matthias Bethge. A sober look at progress in language model reasoning: Pitfalls and paths to
631 reproducibility, 2025. URL <https://arxiv.org/abs/2504.07086>.
- 632
- 633 Yiming Huang, Xiao Liu, Yeyun Gong, Zhibin Gou, Yelong Shen, Nan Duan, and Weizhu Chen. Key-
634 point-driven data synthesis with its enhancement on mathematical reasoning. In *Proceedings of
635 the AAAI 2025*, 2025. URL <https://ojs.aaai.org/index.php/AAAI/article/view/34593>.
- 636
- 637 Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun
638 Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei
639 Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng
640 Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report, 2024. URL <https://arxiv.org/abs/2409.12186>.
- 641
- 642 Shima Imani, Liang Du, and Harsh Shrivastava. Mathprompter: Mathematical reasoning using large
643 language models. In *Proceedings of ACL 2023*, pp. 37–42, 2023. URL [https://aclanthology.
644 org/2023.acl-industry.4.pdf](https://aclanthology.org/2023.acl-industry.4.pdf).
- 645
- 646 inclusionAI. Areal: Ant reasoning reinforcement learning for llms, March 2025. URL <https://github.com/inclusionAI/AReal>.
- 647
- 648 Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec
649 Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint
650 arXiv:2412.16720*, 2024. URL <https://arxiv.org/abs/2412.16720>.

- 648 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
649 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
650 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating
651 Systems Principles*, 2023. URL <https://arxiv.org/abs/2309.06180>.
- 652 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay
653 Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam
654 Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with language
655 models. In *Proceedings of NeurIPS 2022*, 2022. URL [http://papers.nips.cc/paper_files/
656 paper/2022/hash/18abbeef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/18abbeef8cfe9203fdf9053c9c4fe191-Abstract-Conference.html).
- 657 Chen Li, Weiqi Wang, Jingcheng Hu, Yixuan Wei, Nanning Zheng, Han Hu, Zheng Zhang, and
658 Houwen Peng. Common 7b language models already possess strong math capabilities, 2024a.
659 URL <https://arxiv.org/abs/2403.04706>.
- 660 Chengpeng Li, Zheng Yuan, Hongyi Yuan, Guanting Dong, Keming Lu, Jiancan Wu, Chuanqi
661 Tan, Xiang Wang, and Chang Zhou. Mugglemath: Assessing the impact of query and response
662 augmentation on math reasoning. In *Proceedings of ACL 2024*, pp. 10230–10258, 2024b. URL
663 <https://aclanthology.org/2024.acl-long.551.pdf>.
- 664 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif Rasul,
665 Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in ai4maths
666 with 860k pairs of competition math problems and solutions. *Hugging Face repository*, 13:9, 2024c.
667 URL http://faculty.bicmr.pku.edu.cn/~dongbin/Publications/numina_dataset.pdf.
- 668 Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason Weston, and
669 Mike Lewis. Self-alignment with instruction backtranslation. In *Proceedings of ICLR 2024*, 2024d.
670 URL <https://openreview.net/forum?id=1oiJHBRsT>.
- 671 Xuefeng Li, Haoyang Zou, and Pengfei Liu. LIMR: Less is More for RL Scaling. *arXiv preprint
672 arXiv:2502.11886*, 2025. URL <http://arxiv.org/abs/2502.11886>.
- 673 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike,
674 John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *Proceedings of ICLR
675 2024*, 2024. URL <https://openreview.net/forum?id=v8L0pN6E0i>.
- 676 Haoxiong Liu, Yifan Zhang, Yifan Luo, and Andrew C Yao. Augmenting math word problems
677 via iterative question composing. In *ICLR 2024 Workshop on Navigating and Addressing Data
678 Problems for Foundation Models*, 2024. URL <https://arxiv.org/abs/2401.09003>.
- 679 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
680 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
681 approach, 2019. URL <https://arxiv.org/abs/1907.11692>.
- 682 Zichen Liu, Changyu Chen, Wenjun Li, Tianyu Pang, Chao Du, and Min Lin. There may not be
683 aha moment in rl-zero-like training — a pilot study. <https://oatllm.notion.site/oat-zero,>
684 2025a. Notion Blog.
- 685 Zihan Liu, Yang Chen, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Acemath: Advancing
686 frontier math reasoning with post-training and reward modeling, 2025b. URL [https://arxiv.
687 org/abs/2412.15084](https://arxiv.org/abs/2412.15084).
- 688 AI @ Meta Llama Team. The llama 3 herd of models, 2024. URL [https://arxiv.org/abs/2407.
689 21783](https://arxiv.org/abs/2407.21783).
- 690 Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng,
691 Qingwei Lin, Shifeng Chen, Yansong Tang, and Dongmei Zhang. Wizardmath: Empowering
692 mathematical reasoning for large language models via reinforced evol-instruct. In *Proceedings of
693 ICLR 2025*, 2025a. URL <https://openreview.net/forum?id=mMPMHw0d0y>.
- 694 Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Tang, Manan Roongta, Colin Cai,
695 Jeffrey Luo, Tianjun Zhang, Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing
696 o1-preview with a 1.5b model by scaling rl, 2025b. Notion Blog.

- 702 Lingrui Mei, Shenghua Liu, Yiwei Wang, Baolong Bi, Yuyao Ge, Jun Wan, Yurong Wu, and Xueqi
703 Cheng. a1: Steep test-time scaling law via environment augmented generation. *arXiv preprint*
704 *arXiv:2504.14597*, 2025.
- 705
706 Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwu Hu, Yiru Tang, Jiapeng Wang,
707 Xiaoxue Cheng, Huatong Song, Wayne Xin Zhao, Zheng Liu, Zhongyuan Wang, and Ji-Rong Wen.
708 Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems,
709 2024. URL <https://arxiv.org/abs/2412.09413>.
- 710 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke
711 Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time
712 scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.
- 713 OpenAI. OpenAI o3-mini System Card, January 2025. URL [https://cdn.openai.com/
714 o3-mini-system-card-feb10.pdf](https://cdn.openai.com/o3-mini-system-card-feb10.pdf).
- 715
716 Zhuoshi Pan, Yu Li, Honglin Lin, Qizhi Pei, Zinan Tang, Wei Wu, Chenlin Ming, H. Vicky Zhao,
717 Conghui He, and Lijun Wu. Lemma: Learning from errors for mathematical advancement in llms,
718 2025. URL <https://arxiv.org/abs/2503.17439>.
- 719 Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar,
720 and Jason Weston. Iterative reasoning preference optimization. In *Proceedings of*
721 *NeurIPS 2024*, 2024. URL [http://papers.nips.cc/paper_files/paper/2024/hash/
722 d37c9ad425fe5b65304d500c6edcba00-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/d37c9ad425fe5b65304d500c6edcba00-Abstract-Conference.html).
- 723
724 Keiran Paster, Marco Dos Santos, Zhangir Azerbayev, and Jimmy Ba. Openwebmath: An open
725 dataset of high-quality mathematical web text. In *Proceedings of ICLR 2024*, 2023. URL
726 <https://openreview.net/forum?id=jKHmjlpViu>.
- 727 Qwen. Qwq-32b: Embracing the power of reinforcement learning, 2025. URL [https://qwenlm.
728 github.io/blog/qwq-32b/](https://qwenlm.github.io/blog/qwq-32b/).
- 729
730 Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
731 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,
732 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin
733 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi
734 Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,
735 Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL
736 <https://arxiv.org/abs/2412.15115>.
- 737 Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric
738 Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer:
739 Language models can teach themselves to use tools. In *Proceedings of Neurips*
740 *2024*, volume 36, 2024. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
741 d842425e4bf79ba039352da0f658a906-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/d842425e4bf79ba039352da0f658a906-Abstract-Conference.html).
- 742 Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set
743 approach. In *Proceedings of ICLR 2018*, 2018. URL [https://openreview.net/forum?id=
744 H1aIuk-RW](https://openreview.net/forum?id=H1aIuk-RW).
- 745
746 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li, Yu Wu,
747 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language
748 models. *preprint, arXiv:2402.03300*, 2024. URL <https://arxiv.org/abs/2402.03300>.
- 749 Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and
750 Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL
751 <https://arxiv.org/abs/2303.11366>.
- 752 Zhengyang Tang, Xingxing Zhang, Benyou Wang, and Furu Wei. Mathscale: Scaling instruc-
753 tion tuning for mathematical reasoning. In *Proceedings of ICML 2024*, 2024. URL [https://
754 openreview.net/forum?id=Kjww7ZN47M](https://openreview.net/forum?id=Kjww7ZN47M).
- 755
756 Team. Open Thoughts, January 2025. URL <https://open-thoughts.ai>.

- 756 WizardLM Teams. Wizardlm 2, 2024. URL <https://wizardlm.github.io/WizardLM2/>.
- 757
- 758 Shubham Toshniwal, Wei Du, Ivan Moshkov, Branislav Kisanin, Alexan Ayrapetyan, and Igor
759 Gitman. Openmathinstruct-2: Accelerating ai for math with massive open-source instruction data.
760 In *Proceedings of ICLR 2025*, 2025. URL <https://openreview.net/forum?id=mTCbq2QssD>.
- 761 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and
762 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. In
763 *Proceedings of the ACL 2023*, 2023. URL <https://aclanthology.org/2023.acl-long.754/>.
- 764
- 765 Zengzhi Wang, Xuefeng Li, Rui Xia, and Pengfei Liu. Mathpile: A billion-
766 token-scale pretraining corpus for math. In *Proceedings of Neurips 2024*,
767 2024. URL [https://proceedings.neurips.cc/paper_files/paper/2024/hash/
2d0be3cd5173c10b6ec075d1c393a13d-Abstract-Datasets-and-Benchmarks-Track.html](https://proceedings.neurips.cc/paper_files/paper/2024/hash/2d0be3cd5173c10b6ec075d1c393a13d-Abstract-Datasets-and-Benchmarks-Track.html).
- 768
- 769 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
770 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings
771 of Neurips 2022*, pp. 24824–24837, 2022. URL <https://arxiv.org/abs/2201.11903>.
- 772
- 773 Liang Wen, Yunke Cai, Fenrui Xiao, Xin He, Qi An, Zhenyu Duan, Yimin Du, Junchen Liu, Lifu
774 Tang, Xiaowei Lv, Haosheng Zou, Yongchao Deng, Shousheng Jia, and Xiangzheng Zhang.
775 Light-rl: Curriculum sft, dpo and rl for long cot from scratch and beyond, 2025. URL <https://arxiv.org/abs/2503.10460>.
- 776
- 777 Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. Inference scaling laws:
778 An empirical analysis of compute-optimal inference for problem-solving with language models.
779 *preprint, arXiv:2408.00724*, 2024. URL <https://arxiv.org/abs/2408.00724>.
- 780
- 781 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and
782 Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. In
783 *Proceedings of ICLR 2024*, 2024. URL <https://arxiv.org/abs/2304.12244>.
- 784
- 785 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu,
786 Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu,
787 Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert
788 model via self-improvement, 2024. URL <https://arxiv.org/abs/2409.12122>.
- 789
- 790 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok,
791 Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical
792 questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023. URL <https://arxiv.org/abs/2309.12284>.
- 793
- 794 Qianjin Yu, Keyu Wu, Zihan Chen, Chushu Zhang, Manlin Mei, Lingjun Huang, Fang Tan, Yongsheng
795 Du, Kunlin Liu, and Yurui Zhu. Rethinking the generation of high-quality cot data from the
796 perspective of llm-adaptive question difficulty grading, 2025. URL [https://arxiv.org/abs/
2504.11919](https://arxiv.org/abs/2504.11919).
- 797
- 798 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen.
799 Mammoth: Building math generalist models through hybrid instruction tuning, 2023. URL
800 <https://arxiv.org/abs/2309.05653>.
- 801
- 802 Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhao Chen. Mammoth2: Scaling instructions from the
803 web, 2024. URL <https://arxiv.org/abs/2405.03548>.
- 804
- 805 Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. SimpleRL-
806 Zoo: Investigating and Taming Zero Reinforcement Learning for Open Base Models in the Wild.
807 *arXiv preprint arXiv:2503.18892*, 2025. URL <https://arxiv.org/abs/2503.18892>.
- 808
- 809 Yan Zeng, Hanbo Zhang, Jiani Zheng, Jiangnan Xia, Guoqiang Wei, Yang Wei, Yuchen Zhang, Tao
Kong, and Ruihua Song. What matters in training a gpt4-style language model with multimodal
inputs? In *Proceedings of the 2024 Conference of the North American Chapter of the Association
for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp.
7930–7957, 2024.

810 Xiaojiang Zhang, Jinghui Wang, Zifei Cheng, Wenhao Zhuang, Zheng Lin, Minglei Zhang, Shaojie
811 Wang, Yinghan Cui, Chao Wang, Junyi Peng, Shimiao Jiang, Shiqi Kuang, Shouyu Yin, Chaohang
812 Wen, Haotian Zhang, Bin Chen, and Bing Yu. Srpo: A cross-domain implementation of large-scale
813 reinforcement learning on llm, 2025. URL <https://arxiv.org/abs/2504.14286>.

814 Xueliang Zhao, Xinting Huang, Wei Bi, and Lingpeng Kong. Segoo: Sequential subgoal op-
815 timization for mathematical problem-solving. In *Proceedings of ACL 2024*, 2024. URL
816 <https://aclanthology.org/2024.acl-long.407.pdf>.

817 Xueliang Zhao, Wei Wu, Jian Guan, and Lingpeng Kong. Promptcot: Synthesizing olympiad-level
818 problems for mathematical reasoning in large language models, 2025. URL <https://arxiv.org/abs/2503.02324>.

819 Huichi Zhou, Kin-Hei Lee, Zhonghao Zhan, Yue Chen, Zhenhao Li, Zhaoyang Wang, Hamed
820 Haddadi, and Emine Yilmaz. Trustrag: Enhancing robustness and trustworthiness in rag, 2025.
821 URL <https://arxiv.org/abs/2501.00879>.

822 Kun Zhou, Beichen Zhang, Jiapeng Wang, Zhipeng Chen, Wayne Xin Zhao, Jing Sha, Zhichao Sheng,
823 Shijin Wang, and Ji-Rong Wen. Jiuzhang3. 0: Efficiently improving mathematical reasoning
824 by training small data synthesis models. In *Proceedings of Neurips 2024*, 2024. URL <https://neurips.cc/virtual/2024/poster/93252>.

825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Table 7: Evaluation results of 7B committee models.

Models	Base	AIME'24	AIME'25	AMC'23	MATH500	Minerva	Olympiad
<i>Qwen-Based Models</i>							
Qwen2.5-Math-7B-Base Yang et al. (2024)	-	20.7±3.8	8.7±3.9	56.2±5.7	64.3±0.5	17.3±1.9	29.0±0.5
Qwen2.5-Math-7B-Instruct (Yang et al., 2024)	Qwen2.5-Math-7B	15.7±3.9	10.7±3.8	67.0±3.9	82.9±0.1	35.0±0.6	41.3±0.9
Qwen2.5-7B-Instruct (Qwen et al., 2025)	Qwen2.5-7B	12.3±3.2	7.3±3.4	52.8±4.8	77.1±1.2	34.9±1.0	38.7±1.0
s1.1-7B (Muennighoff et al., 2025)	Qwen2.5-7B	19.0±3.2	21.0±5.5	59.5±3.7	80.8±0.6	37.5±1.1	48.2±1.4
Eurus-2-7B-PRIME (Cui et al., 2025)	Qwen2.5-7B	17.8±2.2	14.0±1.7	63.0±3.9	80.1±0.1	37.5±1.0	43.9±0.3
WarriorMath-Qwen-7B-Small	Qwen2.5-7B	24.4±3.3	22.3±5.2	70.6±3.9	81.6±0.5	37.9±3.8	46.7±1.3

A DISCUSSION

A.1 DIFFERENCES FROM OTHER DATA SYNTHESIS STRATEGIES

WarriorMath diverges from conventional knowledge distillation in three fundamental ways:

(1) Data Sources. Traditional distillation typically relies on pre-existing seed datasets (e.g., NuminaMath-1.5 (Li et al., 2024c), DeepScaleR (Luo et al., 2025b)) and then distills solutions from expert models. In contrast, WarriorMath elicits new, challenging problems directly from experts and derives high-quality solutions through competitive interactions among them.

(2) Learning Objectives. While conventional distillation primarily focuses on imitating expert performance, WarriorMath establishes a closed-loop paradigm of *detecting defects and correcting weaknesses*, where the base model’s shortcomings are explicitly identified and iteratively addressed.

(3) Knowledge Integration. Standard distillation usually adopts responses from a single expert, whereas WarriorMath aggregates strengths from multiple experts (e.g., Phi, QwQ) and adaptively balances their influence via an Elo rating mechanism.

A.2 EFFECTIVENESS WITH SIMILAR-SIZED MODELS

To verify that WarriorMath does not rely solely on distillation from larger models, we conducted the experiment using a committee of similar-sized models (7B). We select three models: Qwen2.5-7B-Instruct, Qwen2.5-Math-7B-Instruct, DeepSeek-R1-Distill-Qwen-7B to train the Qwen2.5-7B-Base.

As shown in Table 7, The model trained by the 7B-Committee achieved 24.4% on AIME24, which outperforms the teacher model Qwen2.5-Math-7B-Instruct. This demonstrates that WarriorMath is effective even without larger “teacher” models. Notably, baselines in Table 7 (e.g., s1.1) rely on distillation from stronger teachers (including Gemini and DeepSeek-R1-671B, which outperform our strongest teacher model). Despite this advantage, these baselines still underperform WarriorMath. We attribute this improvement to two key factors: **1.Diversity:** Different models have different knowledge boundaries. An “ensemble” of 7B models covers more ground than a single one. **2.Verification > Generation:** A 7B model can often verify the correctness of a solution (acting as a Judge) or generate a hard problem (acting as an Examiner) that is difficult for itself or peers to solve directly. This allows the system to bootstrap improvement rather than just distilling knowledge.

A.3 COMPUTATIONAL COST AND EFFICIENCY

We now report the full computational profile: **Committee assessment:** 8×A100 GPUs for 96 hours, producing 240M tokens (200K samples). **Training:** 8×A100 GPUs for 160 GPU hours (batch size 512, 448 steps). This is ~30% more efficient than baselines such as LIMR, which required over 3,000 hours.

A.4 CLARIFYING BASELINE FAIRNESS

In Table 1, baselines such as s1.1 (Muennighoff et al., 2025), Light-R1 (Wen et al., 2025), and LIMR (Li et al., 2025) rely on distillation from stronger teachers, including Gemini and DeepSeek-R1 671B (DeepSeek-AI et al., 2025), which surpass our strongest teacher (DeepSeek-R1-Distill-Llama-70B). Despite this advantage, these baselines underperform WarriorMath. For instance, on AIME-2024, WarriorMath-Qwen-7B achieves 48.3% Pass@1, compared to LIMR’s 30.7% and s1.1’s

19.0%. This gap underscores the superiority of WarriorMath’s synthesis paradigm over conventional distillation.

In Table 2, we ensure fair comparisons across data generation methods. For approaches without released datasets (e.g., Evol-Instruct (Luo et al., 2025a), KPDDS (Huang et al., 2025)), we replicate their setups using Llama-3.1-70B-Instruct (Llama Team, 2024) to generate problems at scale matched to WarriorMath. For NuminaMath (Li et al., 2024c) and OpenMathInstruct (Toshniwal et al., 2025), we directly adopt their released sets. To standardize solution generation, we uniformly employ Qwen2.5-Math-72B-Instruct (Qwen et al., 2025) as the solver across all baselines.

A.5 ON THE RISK OF DATA LEAKAGE

All WarriorMath problems are synthesized by expert models through adversarial prompting, rather than sampled from existing datasets. If leakage were to occur, it would imply that the expert models themselves had memorized training data. However, the experts employed (e.g., DeepSeek-R1, Qwen2.5-Math) explicitly document rigorous data curation efforts to mitigate leakage risks (DeepSeek-AI et al., 2025; Qwen et al., 2025).

To further ensure novelty, we apply ROUGE-based semantic similarity screening (Figure 3) to reduce overlap. Consequently, WarriorMath instructions provide diverse, independent training examples that enhance generalization.

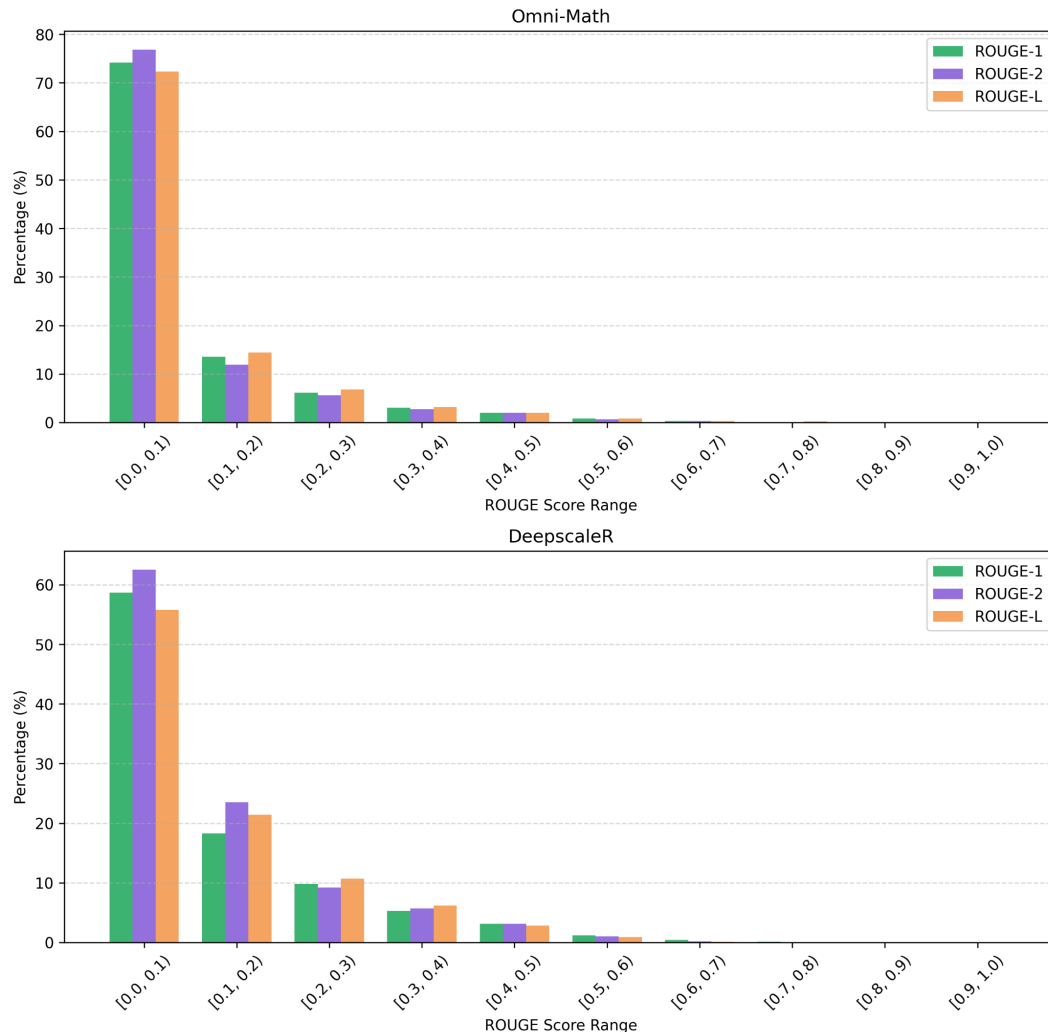


Figure 3: The overlapping rate between the mined instructions and existing training datasets.

B INSTRUCTION DATA SYNTHESIS

B.1 PROMPTS FOR INSTRUCTION MINING

Instruction Mining Prompt

Prompt:

Please act as a world-class expert in designing extremely challenging and diverse math problems. Your goal is to create problems that thoroughly test a model’s reasoning abilities by inducing a variety of potential failure modes (e.g., reasoning, understanding, calculation, or strategy errors).

For each problem you design, please ensure the following:

1. **Quality:** Questions must be well-formatted, clearly structured, and unambiguous.
2. **Difficulty:** Problems should require deep mathematical reasoning and not be solvable via simple pattern recognition or surface-level heuristics.
3. **Diversity:** Problems must span a wide range of mathematical domains, such as algebra, calculus, discrete math, geometry, and others.
4. **Challenge:** Each problem should be adversarially constructed to trigger potential weaknesses in advanced models, such as subtle traps or misleading intermediate steps.

Always provide the final answer enclosed within `\boxed{}` for clarity.

Figure 4: The prompt for generating challenging and diverse math problems.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

B.2 CASE STUDY OF MINING PROBLEMS

Examples of mining problems

Case #1: Let S be the set of all convex quadrilaterals inscribed in the circle $x^2 + y^2 = 25$ with vertices at points having integer coordinates. Determine the maximum possible area of such a quadrilateral Q , given that the sum of the x -coordinates of its vertices equals the sum of the y -coordinates. Express your answer as a reduced fraction $\frac{m}{n}$, where m and n are coprime positive integers, and find $m + n$.

Case #2: Consider triangle ABC with $AB = 13$, $BC = 14$, and $AC = 15$. Let O be the circumcenter and H the orthocenter. Let the circumradius be R . A circle centered at O with radius $R/2$ intersects the nine-point circle at points P and Q . Find the length of PQ .

Case #3: Consider the function $f(x, y) = x^3 + y^3 - 3xy$. Find the maximum value of $f(x, y)$ on the closed disk $x^2 + y^2 \leq 1$.

Case #4: Consider a sequence of n independent coin flips, where each flip has a probability p of landing heads and a probability $q = 1 - p$ of landing tails. Let X be the random variable representing the number of heads in the sequence. Find the probability that X is even.

Case #5: Let $f(x) = x^3 - 3x + 1$. The polynomial $f(x)$ has three real roots, denoted by α, β, γ . Define the sequence $\{a_n\}$ by $a_1 = \alpha + \beta + \gamma$, $a_2 = \alpha^2 + \beta^2 + \gamma^2$, and for $n \geq 3$, $a_n = \alpha^n + \beta^n + \gamma^n$. Find the value of a_{2023} modulo 3.

Case #6: Consider the set $S = \{1, 2, 3, \dots, 100\}$. A subset A of S is called "sum-free" if there do not exist distinct elements $a, b, c \in A$ such that $a + b = c$. Determine the maximum possible number of elements in a sum-free subset of S .

Case #7: A fair six-sided die is rolled repeatedly until a 6 appears. Let X be the number of rolls required. Define the function $f(n)$ as the probability that X is a multiple of n . Find the value of $f(3)$.

Case #8: Determine the value of a_{100} for the sequence defined by $a_1 = 1$ and $a_{n+1} = a_n + \gcd(a_n, n)$ for $n \geq 1$.

Figure 5: Examples of mining problems.

B.3 PROMPT FOR PROBLEM QUALITY EVALUATION AND DISCARDED PROBLEMS EXAMPLES

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Math Instruction Quality Evaluation Prompt

You are a senior university-level mathematics instructor with extensive expertise in advanced topics such as Algebra, Precalculus, Number Theory, Geometry, and Combinatorics. Your task is to **evaluate the quality of mathematical problem statements** based on their clarity, formatting, conceptual soundness, computational complexity, and contextual relevance. Each problem should be scored on a scale of **1 to 10**, and your output must follow a structured JSON format. **Evaluation Criteria:**

- **Clarity and Completeness:** Is the problem clearly stated without ambiguity? Are all necessary variables, conditions, and constraints defined? Is the mathematical notation properly used and well-structured?
- **Conceptual Soundness and Difficulty:** Does the problem involve meaningful, non-trivial mathematical reasoning or advanced concepts? Does it promote critical thinking and apply appropriate mathematical principles?
- **Computational Complexity:** Does the solution process require more than basic arithmetic or trivial computation? Are there non-obvious calculations, transformations, or logical deductions involved?
- **Contextual Relevance and Verifiability:** Is the problem well-grounded in a practical, educational, or theoretical context? Is the problem solvable or verifiable using existing tools and methods? Avoid problems that are vague, proof-based without criteria, or ill-posed.

Scoring Scale:

- **Excellent (9–10):** The instruction is verifiable, properly formatted, and conceptually sound.
- **Good (6–8):** Minor issues in clarity or formatting. Still verifiable and mathematically valid.
- **Average (3–5):** Noticeable flaws in clarity, completeness, or relevance.
- **Poor (1–2):** Ambiguous, improperly defined, unverifiable, or conceptually flawed.

Your Output Format: Your output must be a **JSON list**, where each element is a dictionary with the following keys:

- **instruction:** The original math problem.
- **score:** An integer from 1 to 10 representing your evaluation.
- **reason:** A detailed explanation justifying the score based on the criteria above.

Figure 6: Prompt used to evaluate the quality of mathematical instructions, including scoring criteria and output format. Only instructions rated 6 or higher are considered suitable for use in further steps.

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

Examples of discarded problems

We provide examples of low-quality problems that were filtered out during problem selection, categorized according to our criteria (Verifiability, Proper Formatting, Clarity):

- **(1) Unverifiable Problem:**
Write a python code that run a math function like "log(base , number)".
Reason: This problem is not a mathematical question with a concrete answer, and cannot be automatically evaluated.
- **(2) Poor Formatting:**
With what polynomial function equation do you want to calculate the vertex of the graph?
 $f(x)=2x^2-x^3+5x^4$
Reason: The input mixes natural language with improperly rendered HTML, leading to parsing and readability issues.
- **(3) Incomplete Problem:**
Chef's portion took 15 seconds, and the assistant's portion took 45 seconds.
Reason: The question is incomplete and lacks a clear task or objective for the model to solve.

Figure 7: Examples of discarded problems during the filtering process.

B.4 VOTE AND ELO RATING

The reviewer B is required to respond to the Examiner’s question, while the Examiner A must also provide an answer to its own instruction. Then we can calculate the *local score* for each response:

$$x_{A>B}^i = \frac{t_A}{t_A + t_B} \quad x_{B>A}^i = \frac{t_B}{t_A + t_B} \quad (8)$$

where $x_{A>B}^i$ and $x_{B>A}^i$ are the local scores for A’s and B’s responses to the instruction i . $x_{A>B}^i$ represents the percentage of votes that candidate A receives, while $x_{B>A}^i$ similarly represents the percentage of votes that candidate B receives. t_A and t_B are the number of votes which A and B win.

However, relying solely on the *local score* to select the winner can be problematic. In some cases, a weaker model may receive more votes than a stronger one, even though its responses are not significantly better. This can occur because the *local score* may not fully capture the quality of the model’s performance, especially in situations where the voting is influenced by factors, such as randomness or bias from LLM judges.

To address this limitation, we propose considering both local contingency and global consistency in the decision-making process. Instead of directly basing our analysis on the immediate voting outcomes, we introduce the concept of the *global score* — specifically, the Elo rating, which provides a more comprehensive reflection of a model’s relative performance over time and across various evaluations. The Elo rating system, originally developed to calculate the relative skill levels of players in two-player games (such as chess), has been successfully adapted to assess the performance of competitors in a range of competitive scenarios, including esports and other skill-based games.

By incorporating the Elo rating, we account for both local performance in individual contests and global performance across multiple rounds, providing a more robust and accurate measure of a model’s overall ability. This helps to mitigate the risk of weak models winning based on isolated, potentially unrepresentative votes:

$$X_{A>B}^{Elo} = \frac{1}{1 + 10^{(R_B - R_A)/400}} \quad (9)$$

$$X_{B>A}^{Elo} = \frac{1}{1 + 10^{(R_A - R_B)/400}}$$

where $X_{A>B}^{Elo}$ and $X_{B>A}^{Elo}$ indicate the expected probabilities of A defeating B and B defeating A, respectively. R_A and R_B are the Elo rating of A and B, which are updated dynamically and iteratively. Given the battle result of A and B on an instruction i , we update them by:

$$R_A \leftarrow R_A + K \times (s_{A>B}^i - X_{A>B}^{Elo})$$

$$R_B \leftarrow R_B + K \times (s_{B>A}^i - X_{B>A}^{Elo}) \quad (10)$$

where $s_{A>B}^i$ and $s_{B>A}^i$ are the actual score of the battle result of player A and B (1 for a win, 0.5 for a draw, and 0 for a loss). The factor K controls the sensitivity of rating changes.

Based on Equation 8 and Equation 9, we can obtain the final score of A’s response for instruction i :

$$e_A^i = \sum_{B \in Com \setminus A} \alpha X_{A>B}^{Elo} + (1 - \alpha) x_{A>B}^i \quad (11)$$

where Com is the set of all the competitors and ‘\’ is the subtraction operation. α is the coefficient to balance the local contingency and global consistency.

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

Pair-wise Answer Quality Evaluation Prompt

Prompt:
Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user prompt displayed below.
You will be given a user prompt, a reference answer, and the assistant’s answer. Your job is to compare the assistant’s answer with the reference one and assign a score.
For each user prompt, carry out the following steps:

1. Consider if the assistant’s answer is helpful, relevant, and concise.
 - **Helpful** means the answer correctly responds to the prompt or follows the instructions.
 - **Relevant** means all parts of the response closely connect or are appropriate to what is being asked.
 - **Concise** means the response is clear and not verbose or excessive.
2. Then consider the creativity and novelty of the assistant’s answer when needed.
3. Identify any missing important information in the assistant’s answer that would be beneficial to include when responding to the user prompt.
4. After providing your explanation, you must rate the assistant’s answer on a scale of 1 to 10, where a higher score reflects higher quality.

Guidelines for Scoring:

- **Assistant’s Answer >> Reference Answer (7–10):** The assistant’s answer is significantly or slightly better than the reference answer.
- **Assistant’s Answer == Reference Answer (5–6):** The quality of assistant’s answer is relatively the same as that of the reference answer.
- **Assistant’s Answer << Reference Answer (1–4):** The assistant’s answer is significantly or slightly worse than the reference answer.

User Prompt:
{instruction}

Reference Answer:
{reference}

Assistant’s Answer:
{response}

Use double square brackets to format your scores, like so: [[7]].

Figure 8: The prompt for the evaluation of pairwise comparison.

1296 C THE USE OF LARGE LANGUAGE MODELS
1297

1298 We used a Large Language Model (LLM) only as a writing assistant to polish the language of the
1299 manuscript (*e.g.*, grammar refinement, style adjustment, and clarity improvement). The research
1300 ideas, methodology design, experiments, and analysis were entirely conceived, implemented, and
1301 validated by the authors without reliance on the LLM. The LLM did not contribute to research
1302 ideation, experimental design, or result interpretation.

1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349