# RETOOL: REINFORCEMENT LEARNING FOR STRATEGIC TOOL USE IN LLMS

#### **Anonymous authors**

000

001

002 003 004

005

006 007 008

010

012

013

014

015

016

018

019

020

021

023

024

029

031

033

040

041

044 045 Paper under double-blind review

# **ABSTRACT**

While reasoning models trained with reinforcement learning (RL) excel in reasoning, they struggle in scenarios requiring structured problem-solving, such as geometric reasoning, concise computation, or complex equation solving—areas where computational tools like code interpreters (CI) demonstrate distinct advantages. To bridge this gap, we propose **ReTool**, which enhances long-form reasoning with tool-integrated learning, including two key features: (1) dynamic interleaving of real-time code execution within natural language reasoning processes, and (2) an automated RL paradigm that allows policy rollouts with multi-turn real-time code execution and teaches the model in learning when and how to invoke tools based on outcome feedback. ReTool employs a systematic training framework, beginning with synthetic code-augmented long-form reasoning data for cold-start training. Subsequent RL training leverages task outcomes as rewards to iteratively refine the model's tool use strategy, enabling autonomous discovery of optimal tool invocation patterns without human priors. Experiments on challenging MATH Olympiad benchmark AIME demonstrate ReTool's superiority: Our 32B model achieves 67% accuracy with 400 training steps, outperforming text-based RL baseline (40% accuracy, 1080 steps) in performance and efficiency. Remarkably, ReTool-32B attains 72.5% accuracy in extended settings, surpassing OpenAI's o1-preview by 27.9%. Further analysis reveals generalization to broader tool-use scenarios and emergent behaviors such as code self-correction, signaling an "aha moment" in which the model autonomously masters adaptive tool use. These findings highlight the promise of outcome-driven tool integration for advancing complex mathematical reasoning and offer new insights into hybrid neuro-symbolic systems.

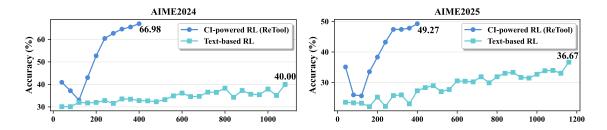


Figure 1: AIME scores of ReTool and text-based RL baseline based on the Qwen2.5-32B-Instruct model.

#### 1 Introduction

Reinforcement learning (RL) has recently become a popular paradigm for enhancing the reasoning capabilities of large language models (LLMs), enabling them to explore and refine long chains of thought (CoT)

(Wei et al., 2022; Yao et al., 2023; Luong et al., 2024; Zhang et al., 2024). Reasoning models such as OpenAI of (OpenAI et al., 2024) and DeepSeek R1 (DeepSeek-AI et al., 2025) demonstrate strong performance in pure text-based reasoning tasks by learning to self-correct and engage in more deliberate, analytical thinking (Claude, 2025; Team, 2025d;a). These advances suggest early signs of metacognitive control, where models not only reason, but also monitor and revise their reasoning process.

Despite these advances, reasoning LLMs equipped with long chains of textual reasoning processes (Ouyang et al., 2022) still show notable limitations in tasks that require precise numerical calculation or symbolic manipulation, such as geometric reasoning, precise computation, or complex equation solving. In contrast, computational tools, such as code interpreters (CI), can empower models with symbolic computation capabilities that go far beyond pure text-based reasoning. Unlike textual CoT (Wei et al., 2023) methods that rely solely on internal language patterns, code interpreters provide a formal and executable interface for enumeration, verification, and precise computation. This not only enables exact numeric validation of intermediate steps—dramatically reducing the ambiguity and compounding error often seen in textual reasoning (Chen et al., 2023; Wang et al., 2023b), but also allows models to expand their solution search space via programmable exploration.

Recent works have explored prompting and supervised fine-tuning methods (Chen et al., 2025b; Pan et al., 2023) to equip LLMs with tool-use capabilities. However, these approaches are limited to imitating the specifically-curated data distribution, often failing to generalize beyond seen patterns or adaptively decide when and how to invoke external tools. As a result, models may misuse tools or fall back on brittle heuristics that are not robust across diverse problem settings. To overcome these limitations, RL offers a principled solution: it enables models to explore flexible reasoning trajectories and learn tool-use strategies guided by outcome-based feedback. This paradigm not only incentivizes correct solutions, but also allows the model to discover nuanced behavioral patterns—such as how to recover from tool execution mistakes via self-correction, decide when to effectively invoke tool execution during the long-chain reasoning process.

In this work, we embrace the RL paradigm and introduce **ReTool**, a **Tool**-augmented **Re**inforcement learning framework explicitly designed to guide LLMs towards optimal strategies for leveraging external computational tools during reasoning. ReTool consists of two key components: First, we develop a data construction pipeline to curate a high-quality cold-start dataset that explicitly demonstrates when and how to invoke the code interpreter. This teaches the model an initial competency in tool usage and execution result analysis. Then, we apply tool-enhanced reinforcement learning to train the model in discovering optimal tool manipulation reasoning strategy and adjusting its behavior through outcome-based rewards, going beyond what can be captured by supervised learning alone. During long-chain reasoning, the policy model rolls out by flexibly writing code blocks and achieving real-time execution results from a sandbox-style code interpreter to assist subsequent thinking.

We evaluate ReTool on the challenging MATH Olympiad benchmarks AIME2024 and AIME2025. Building on Qwen2.5-32B-Instruct (Yang et al., 2024a), our model achieves 67.0% accuracy on AIME2024 with only 400 training steps, significantly outperforming the text-based RL baseline, which achieves 40.0% accuracy with 1080 training steps. These substantial gains highlight that explicitly modeling tool-use as part of the decision process not only pushes limits of model reasoning but also enhances training efficiency. Furthermore, when trained on DeepSeek-R1-Distill-Qwen-32B (DeepSeek-AI et al., 2025), our model demonstrates further improvements, surpassing competitive baselines such as QwQ-32B-Preview (Team, 2025d), s1-32B (Muennighoff et al., 2025), and OpenAI o1-preview (OpenAI, 2024). This suggests that the RL training process inspires more efficient problem-solving strategies. We further conduct a comprehensive analysis, including ablation study, extension to the web search domain, CI cognitive behavior through RL training. This analysis highlights several key findings: our model demonstrates enhanced code utilization capabilities, enabling it to employ more accurate and complex code snippets; It also learns to invoke tools appropriately, select tool adaptively, structure tool calls effectively, and iteratively refine reasoning through emergent code self-correction capabilities. Overall, our main contributions are summarized as follows:

- We propose ReTool, a novel reinforcement learning framework that integrates code interpreter execution into the reasoning loop of LLMs. To equip the model with foundational capabilities for invoking the code interpreter, we curate a high-quality cold-start dataset through our developed pipeline. Furthermore, we design a reinforcement learning framework that supports interleaved code execution during rollout, enabling the model to iteratively explore, refine, and optimize its reasoning strategies through tool-augmented interactions guided by feedback from a sandboxed code interpreter.
- As shown in section 3.6, we conduct comprehensive empirical and behavioral analyses, and observe several key findings: (1) After RL training, the response length is reduced by approximately 40% compared to that prior to training, showcasing the potential reasoning token efficiency of tool-integrated reasoning; (2) Our approach can generalize to broader tool-use scenarios like web search; (3) During RL training, the code ratio, code lines and correct code counts show increase trends, and the code invocation timing becoming shifts earlier, indicating the improved code use capabilities and strategic tool usage development; (4) Emergent behaviors like code self-correction and adaptive tool selection can be observed during RL phase, bringing more advanced tool-augmented reasoning patterns.

# 2 METHODOLOGY

In this section, we introduce ReTool. We begin with an overview of the framework, followed by a description of the cold-start training process, including the data construction pipeline and supervised fine-tuning (section 2.2). We then outline our reinforcement learning pipeline, enhanced by a code interpreter sandbox, to further enhance strategic tool usage development (section 2.3).

#### 2.1 Overview

Our approach consists of two stages: cold-start supervised fine-tuning followed by reinforcement learning with interleaved code execution rollout. Firstly, we collect data through our designed pipeline for cold-start supervised fine-tuning (SFT), which provides a robust initialization for the reinforcement learning (RL) phase. To enhance our model's tool utilization capabilities, we introduce a specialized tool-using RL pipeline that enhances the model's ability to appropriately select and apply tools during the reasoning process.

#### 2.2 COLD-START FOR TOOL-INTEGRATED REASONING FOUNDATION

We designed a pipeline to collect and curate high-quality data. Specifically, we begin by gathering existing mathematical reasoning data from diverse sources, including open-source datasets such as Open-Thoughts (Team, 2025c). Subsequently, we implement a dual-verification approach combining human expert curation and Deepseek-R1 (DeepSeek-AI et al., 2025) evaluation to filter invalid data. Through these steps, we collect a high-quality text-based reasoning dataset, denoted as  $\mathcal{D}_{init}$ .

Based on  $\mathcal{D}_{init}$ , we further construct code-integrated reasoning data automatically. We first utilize a structured prompt template (detailed in Appendix Figure 9) for transformation, which modifies the original thinking process by replacing manual calculation steps that can benefit from code execution with the corresponding code snippets and their interpreter's execution results. Following this initial transformation, we apply a two-stage verification protocol. The first stage focuses on format verification, which improves readability and ensures consistent syntax that that enables the efficient detection of computational tool invocation triggers during subsequent reinforcement learning phases. The second stage entails answer verification, where we eliminate data samples whose final outputs do not align with the correct solutions to the mathematical problems. Finally, we collect a dataset  $\mathcal{D}_{CI}$  that consist of code-augmented long-form reasoning traces.

ReTool employs supervised fine-tuning to learn when and how to invoke the code interpreter from the aforementioned dataset  $\mathcal{D}_{CI}$ , thereby enhancing model's capability to appropriately utilize computational tools.

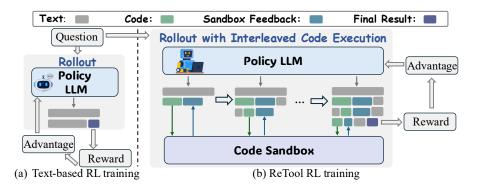


Figure 2: Demonstration of text-based RL training process and ReTool's RL training process.

#### 2.3 RETOOL: REINFORCEMENT LEARNING FOR STRATEGIC TOOL USE

# 2.3.1 Training Algorithm

We train ReTool based on PPO (Schulman et al., 2017), it updates policy with the following objective:

$$\mathcal{J}_{\text{PPO}}(\theta) = \mathbb{E}_{(q,a) \sim \mathcal{D}, o \leq t} \sim \pi_{\theta_{\text{old}}}(\cdot \mid q)} \left[ \min \left( \frac{\pi_{\theta}(o_t \mid q, o_{< t}; \mathcal{CI})}{\pi_{\theta_{\text{old}}}(o_t \mid q, o_{< t}; \mathcal{CI})} \hat{A}_t, \operatorname{clip} \left( \frac{\pi_{\theta}(o_t \mid q, o_{< t}; \mathcal{CI})}{\pi_{\theta_{\text{old}}}(o_t \mid q, o_{< t}; \mathcal{CI})}, 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_t \right) \right]$$

$$(1)$$

where  $\pi_{\theta}$  is policy model,  $\pi_{\theta_{\text{old}}}$  is reference model,  $\pi_{\theta}(o_t \mid q, o_{< t}; \mathcal{CI})$  represents the rollouts with interleaved code execution and feedback from code interpreter. We modify PPO to better adopt tool integrated reasoning. During training, the policy LLM will collaborate with a code sandbox to generate rollouts with multi-turn real-time code execution for solving given problems. We implement a rule-based outcome reward to enable the model with the flexibility to autonomously explore and develop strategies for code usage awareness, code selection, timing of code invocation, and further diverse behaviors.

**Reward Design** To teach the model in learning when and how to invoke tools, we implement a rule-based accuracy reward to optimize the model. The accuracy reward evaluates response correctness. We require the model to present final answers in a specified format (e.g., within \boxed{}), enabling reliable rule-based verification. The reward is formulated as:

$$R(a, \hat{a}) = \begin{cases} 1, & \text{is\_equivalent}(a, \hat{a}) \\ -1, & \text{otherwise} \end{cases}$$
 (2)

where a and  $\hat{a}$  represent the ground-truth answer and the predicted answer, respectively. We simplify the reward design aim to alleviate reward hacking and promote more diverse problem-solving behaviors based on mere outcome feedback without considering code executability reward.

**Rollout with Interleaved Code Execution** To facilitate the integration of reasoning and executable code within the model, we propose a rollout approach that dynamically supports interleaved real-time code execution with natural language reasoning processes. As depicted in Figure 2 (b), our rollout process differs from the conventional approach, which typically generates only text-based reasoning (as shown in Figure 2 (a)). By contrast, our rollout approach integrates the collaboration of a policy LLM with an external code sandbox, enabling the production of hybrid content that combines text, code snippets, and real-time interpreter feedback. Concretely, we utilize a prompt template (Appendix Figure 8) to guide the model in interacting with the code sandbox by utilizing tags <code></code> to explicitly mark the boundaries of generated codes. During the rollout process, policy model generate text-based reasoning  $t_1$  when a code termination trigger (</code>) is detected, the generation pause and the generated code  $c_1$  is parsed and

send to code sandbox environment for execution. Upon completion, the sandbox's output  $f_1$  (successful results or error messages) is filled within <interpreter></interpreter> tags and fed back to the model, which continues generating the rollout until either providing a final answer o or producing a new code snippet, ultimately producing a hybrid reasoning trajectory  $[t_1 \oplus c_1 \oplus f_1 \oplus ... \oplus o]$ .

Notably, our approach returns both successful code execution results and interpreter error messages to the model. This dynamic feedback mechanism enables the model to iteratively explore, refine, and optimize its reasoning and tool usage strategies.

# 2.3.2 Training Details

**Cold-start & RL** For training, we employ the VeRL framework. We adopt PPO as our RL method. We train our model on curated cold-start data for two epochs. Regarding hyperparameters, we utilize the AdamW optimizer with an initial learning rate of 1e-6. We define the expected maximum sequence length as 16384 tokens. For training, the mini-batch size is set to 512, and the KL coefficient is set to 0.0. We use Qwen2.5-32B-Instruct (Qwen et al., 2025) as the main backbone. All experiments are conducted on NVIDIA H20 GPUs.

**Interpreter Feedback Mask.** We mask out the <interpreter></interpreter> feedback output from the loss computation. This sandbox-based output masking approach blocks external tokens from interfering with loss calculations, ensuring training stability and preserving the model's inherently generated coherent reasoning sequences from disruption.

**KV-Cache Reuse.** In order to reduce the memory cost during rollout, when each time the code termination trigger (</code>) is detected, we will cache all the KV-cache before code execution and only calculate and append the KV-cache from the interpreter feedback (<interpreter></interpreter>). This will largely reduce the KV-cache for each rollout.

**Sandbox Construction.** To accelerate the RL training process, we design a asynchornous code sandbox environment. The sandbox pods function as workers in a pool, independently pulling tasks based on their current capacity, creating an efficient load-balancing mechanism. This distributed asynchronous approach accelerates RL training by enabling parallel environment interactions across multiple threads, It prevents slower threads from creating bottlenecks and ensures optimal resource utilization, maintaining continuous throughput during the training process.

## 3 EXPERIMENT

#### 3.1 EVALUATION SETUP

To ensure a stable evaluation, we repeat the evaluation set AIME2024 and AIME2025 32 times, GPQA (Diamond) (Rein et al., 2023) 8 times, MATH500 (Hendrycks et al., 2021) 4 times, GSM8K (Cobbe et al., 2021) 2 times, and report the overall average accuracy to estimate pass@1. The inference hyperparameters of evaluation are set to temperature 1.0 and top-p 0.7. We compare ReTool with competitive baselines, including Qwen2.5-Math-72B-Instruct (Yang et al., 2024b), Qwen2.5-Math-72B-Instruct-TIR (Yang et al., 2024b), Sky-T1 (Team, 2025b), DeepSeek-R1-Zero-Qwen-32B (DeepSeek-AI et al., 2025), QwQ-32B-Preview (Team, 2025d), s1-32B (Muennighoff et al., 2025), OpenAI o1-preview (OpenAI, 2024). To verify the effectiveness of our ReTool, we also compare the performance with RL without tool-using, i.e. Text-based RL (Qwen2.5-32B-Instruct). And for the results of baselines, we report the avg@k by coping from corresponding literature source as pass@1.

Table 1: Main results on benchmarks in mathematic and STEM domains.

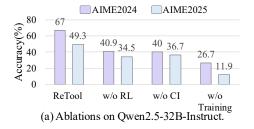
| Model                                 | AIME2024 | AIME2025 | GSM8K | MATH500 | GPQA |
|---------------------------------------|----------|----------|-------|---------|------|
| Qwen2.5-Math-72B-Instruct             | 30.0     | -        | 95.9  | 85.9    | -    |
| Qwen2.5-Math-72B-Instruct-TIR         | 40.0     | -        | 95.8  | 88.1    | -    |
| Sky-T1                                | 43.3     | -        | -     | 82.4    | 56.8 |
| OpenAI o1-preview                     | 44.6     | 37.9     | -     | 85.5    | 73.3 |
| DeepSeek-R1-Zero-Qwen-32B             | 47.0     | -        | -     | 94.3    | 62.1 |
| QWQ-32B-Preview                       | 50.0     | 33.5     | -     | 90.6    | 54.5 |
| s1-32B                                | 56.7     | -        | -     | 93.0    | 59.6 |
| ReTool (Qwen2.5-32B-Instruct)         | 67.0     | 49.3     | 95.9  | 93.1    | 58.7 |
| ReTool (DeepSeek-R1-Distill-Qwen-32B) | 72.5     | 54.3     | 96.3  | 95.2    | 62.3 |

#### MAIN RESULTS

As shown in Table 1, ReTool enables the LLM to flexibly leverage the code interpreter during the RL stage, leading to substantial performance improvements. Specifically, ReTool (Qwen2.5-32B-Instruct) achieves accuracies of 67.0% on AIME2024, 49.3% on AIME2025, 95.9% on GSM8K, 93.1% on MATH, and 58.7% on GPQA with only 400 training steps. This markedly outperforms most strong baselines with much larger parameter sizes, such as Qwen2.5-Math-72B. These findings indicate that the tool-integrated learning paradigm employed by ReTool enhances the model's reasoning capabilities significantly. Furthermore, on AIME2024, ReTool (Qwen2.5-32B-Instruct) surpasses the competitive baseline s1-32B by 10.3%. Similarly, on AIME2025, it achieves an 11.4% gain over OpenAI's o1-preview. When combined with a more advanced backbone, ReTool (DeepSeek-R1-Distill-Qwen-32B) further improves performance, achieving scores of 72.5% on AIME2024 and 54.3% on AIME2025. These results suggest that more effective problem-solving strategies are discovered during the RL training process.

# 3.3 ABLATION STUDY

To further assess the effectiveness of ReTool, we conduct an ablation study and compare against several variants built on Qwen2.5-32B-Instruct: (1) w/o RL: the cold-start model that still incorporates the code interpreter; (2) w/o CI: a text-based RL method initialized with a text-only cold-start SFT to ensure fairness; (3) w/o Training: the original base model. As illustrated in Figure 3 (a), removing either the RL stage (w/o RL) or the CI integration (w/o CI) leads to a notable drop in performance on AIME2024&2025. Furthermore, our cold-start model achieves 40.9% accuracy on AIME2024, which is comparable to the text-based RL baseline (40.0%) and significantly higher than the base model (26.7%). These findings demonstrate that our curated dataset successfully captures tool-usage patterns within executable reasoning traces, and that CI-integrated training further enhances reasoning performance.



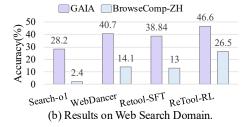


Figure 3: (a) Ablation study; (b) Results on the web search domain with same backbone on Qwen2.5-32B.

# 3.4 EXPANSION TO WEB SEARCH SCENARIO

To assess whether our strategy generalizes to broader tool-use scenarios, we further conducted experiments in the Web Search domain. Specifically, we applied ReTool with the Bing Search tool (following the MCP tool-use definition) on two widely used benchmarks: GAIA (Mialon et al., 2023) and BrowseComp-ZH (Zhou et al., 2025). As illustrated in Figure 3(b), when using the same backbone (Qwen2.5-32B), ReTool consistently outperforms competitive baselines built on the same backbone, including WebDancer (Wu et al., 2025) and Search-o1 (Li et al., 2025a). These findings highlight ReTool's generalizability and effectiveness beyond mathematical reasoning tasks, underscoring its adaptability across diverse tool-use environments.

# 3.5 "AHA MOMENT" OF CODE SELF-CORRECTION.

Interestingly, our model exhibits an emergent ability to self-correct non-executable code, despite the absence of explicit training data for code self-correction. As shown in Figure 6 in Appendix, the model initially produced code that failed to execute due to the undefined function "greedy()". Upon receiving feedback from the interpreter, the model recognized the error and responded with the reflection: "Oops, the functions need to be defined in the same scope. Let's correct that." It then proceeded to generate a revised, executable version of the code that included all necessary function definitions. This emergent behavior suggests that reinforcement learning can foster metacognitive capabilities, enabling the model to iteratively refine its generated code to address more complex problems.

# 3.6 COGNITIVE ANALYSIS

We present a comprehensive analysis, including the dynamics of code interpreter (CI)-related behaviors throughout the RL and the differences in code purpose before and after RL.

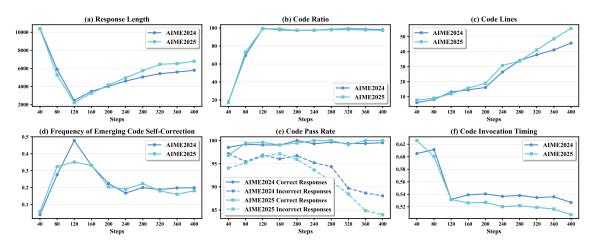


Figure 4: CI-related behavior evolution during RL training.

**CI-related Behavior Evolution.** To gain deeper insights into the RL process of ReTool, we systematically evaluated CI-related metrics. Specifically, we computed various metrics by analyzing model-generated outputs on the AIME2024 and AIME2025 datasets based on each saved checkpoint during RL training:

• Response Length (Figure 4 (a)): We calculated the average response length and observed a distinct trend: the response length initially declines sharply, later followed by a relatively gentle increase. We attribute

the initial decline to the replacement of complex computational processes with more concise code, while the subsequent rise is likely due to the emergence of more diverse and complex code behaviors during RL training. Notably, the final average response length remains 40% shorter than that before RL training (i.e., from 10k to 6k). This suggests that the CI-powered reasoning approach potentially enhances efficiency of reasoning token utilization ratio by replacing intricate computational processes with code.

- Code Ratio (Figure 4 (b)): The ratio of responses that contain code are also calculated. Analysis reveals that throughout the RL training process, the average code ratios exhibit a total upward trend and end with covering nearly 98% percent of all questions. This suggests that the model's proficiency in code utilization improved progressively during the RL process, facilitating strategic tool usage development.
- Code Lines (Figure 4 (c)): The lines of generated code reflects its complexity to some extent. Observations show that the average code lines in responses exhibits a consistent upward trend throughout training. By the end of RL training, the final average code lines is nearly fivefold higher than that before RL training. This trend suggests that the model has learned more complex code strategies during the RL phase.
- Frequency of Emerging Code Self-Correction (Figure 4 (d)): We approximate the frequency of code self-correction by detecting specific turning cues (e.g., oops, wait, or correcting) that appear between a failed code block and its immediate successor. As shown, the frequency of self-correction emerges early and peaks during the initial stages of training, then gradually decreases and stabilizes. This trend suggests that when the model's code generation capability is still limited, it learns and leverages self-correct more frequently. As RL training, improves in producing correct code directly, thereby reducing the necessity for correction. This phenomenon reflects an emergent behavior: self-correction is prevalent at the outset but becomes less necessary as the model gains competence.
- Code Pass Rate (Figure 4 (e)): The CI-powered reasoning process involves generating intermediate code that may initially be incorrect, followed by iterative refinement based on interpreter feedback to produce executable code, so we report the average pass rate of last code in incorrect responses. Our analysis reveals that the code pass rate for correct responses remains consistently high, approaching 100%, while the code pass rate for incorrect responses exhibits a declining trend. This pattern suggests that code executability impacts the reasoning process and final result.
- Code Invocation Timing (Figure 4 (f)): We also calculate the code invocation timing, which is determined by dividing the start position of code by the total length of the response. This metric reflects the timing of code invocation within the response. The results show that the code invocation timing advances during the RL training process, indicating that the model learns to determine the timing for tool usage.

**Code Purpose Analysis.** We also analysis the differences in code purposes before and after RL training, which reflects the types of code. We employ Claude4-Sonnet (Anthropic, 2025) to classify the primary purpose of code snippets based on their contextual information, then compute the frequency of code purposes that appear more than once, and the results are depicted in Figure 5. The word clouds reveal that calculation

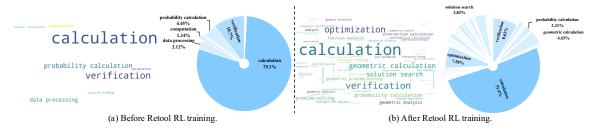


Figure 5: Code purpose analysis.

and verification are the dominant purposes of code in CI-powered reasoning. After RL training, the code purposes in our model become more diverse, which demonstrates the metacognitive development of adaptive tool selection and enhances the generalizability of ReTool to a broader range of problems.

# 4 BACKGROUND AND RELATED WORK

#### 4.1 LLM REASONING

Recent advancements in large language models (LLMs) (Wei et al., 2022; Yao et al., 2023; Luong et al., 2024; OpenAI et al., 2024; Team, 2025a; DeepSeek-AI et al., 2025; xAI, 2023; Claude, 2025; Team et al., 2023; Yang et al., 2024a) indicate significant progress toward cognitive abilities similar to human metacognition through Chain-of-Thought (CoT) prompting, which enhances the reasoning capabilities of LLMs by leveraging step-by-step natural language reasoning Wei et al. (2023); Wang et al. (2023a); Wang & Zhou (2024); Huang et al. (2025); Chen et al. (2025a). Building upon this foundation, recent research has shifted focus from train-time scaling to test-time scaling (Snell et al., 2024), where additional computational resources are allocated during inference to enable the generation of intermediate reasoning steps. Techniques such as stepwise preference optimization (Lai et al., 2024), Monte Carlo Tree Search (MCTS) (Xie et al., 2024), and reinforcement learning (Luong et al., 2024; DeepSeek-AI et al., 2025) have been employed to improve multi-step and long-form mathematical reasoning. Complementing CoT, Program-of-Thought (PoT) reasoning, introduced by Chen et al. (2023) and Gao et al. (2023), integrates external computational tools to simplify and validate complex reasoning steps, resulting in enhanced accuracy.

# 4.2 TOOL INTEGRATED REASONING

Tool-integrated reasoning was first introduced to help LLMs solve computationally intensive mathematical problems with the integration of programming strategies (Chen et al., 2023; Yue et al., 2023; Jin et al., 2025; Song et al., 2025; Wang et al., 2024). Building on this line of work, Wang et al. (2023b) introduced an iterative framework that integrates textual reasoning with code execution, while Chen et al. (2025b) further advanced this direction by applying supervised fine-tuning on self-curated code-integrated CoT data. However, this approach is inherently limited by its reliance on the specific data distribution, and cannot learn adaptive strategies for tool use—such as determining when and how to invoke tools—through reinforcement learning (RL). A concurrent work (Li et al., 2025b) applied RL to learn tool usage strategies on Qwen2.5-Math models (Yang et al., 2024b) at 1.5B and 7B scales, but the performance remained suboptimal. We further scale up this line of research and propose ReTool, a framework that leverages RL to strategically determine when and how to invoke the code interpreter. Our method outperforms Qwen-Math-72B-TIR (Yang et al., 2024b) and o1-preview (OpenAI, 2024) significantly on AIME2024 and AIME2025. We also present a comprehensive analysis of the learned tool-use behaviors and highlight several key findings regarding the model's cognitive patterns in code invocation after ReTool training.

# 5 CONCLUSION

In this paper, we propose ReTool, a novel reinforcement learning framework that empowers LLMs to self-enhance their mathematical reasoning capabilities through effective Code Interpreter utilization. Our comprehensive experiments on various benchmarks demonstrate that ReTool not only achieves superior accuracy compared to conventional text-based RL approaches, but also converges with significantly fewer training steps. Further analyses also demonstrate generalization to broader tool-use scenarios and emergent behaviors such as code self-correction. Through careful data curation and our specialized tool-using pipeline, ReTool enables models to develop sophisticated computational intervention strategies, paving the way for more efficient and powerful tool-augmented reasoning in LLMs.

#### ETHICS STATEMENT

This work does not involve human subjects, personal data, or sensitive information. All experiments were conducted on publicly available datasets in mathematics and STEM domains (e.g., AIME, GSM8K, MATH, GPQA). The datasets were used strictly for research purposes and do not contain identifiable private information. Our methods focus on improving the reasoning capabilities of LLMS by integrating reinforcement learning with external computational tools. We affirm compliance with the ICLR Code of Ethics and confirm that no part of this research raises concerns regarding privacy, discrimination, or conflict of interest.

# REPRODUCIBILITY STATEMENT

We have made every effort to ensure the reproducibility of our work. The full methodological details of ReTool, including cold-start data construction, supervised fine-tuning, reinforcement learning setup, reward design, and evaluation protocols, are described in Sections 2 and 3 of the main text. Comprehensive experimental results across multiple datasets and ablation studies are reported in Tables 1–2 and Figures 3–5. Additional training details such as optimizer, learning rate, batch size, sequence length, masking strategy, and sandbox design are provided in Section 2.3.2. These details should enable independent researchers to replicate our experiments. And the training data and the evaluation scripts are available in the supplementary material.

### REFERENCES

Anthropic. Claude sonnet 4. https://www.anthropic.com/claude/sonnet, 2025.

- Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models, 2025a. URL https://arxiv.org/abs/2503.09567.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks, 2023. URL https://arxiv.org/abs/2211.12588.
- Zhipeng Chen, Yingqian Min, Beichen Zhang, Jie Chen, Jinhao Jiang, Daixuan Cheng, Wayne Xin Zhao, Zheng Liu, Xu Miao, Yang Lu, Lei Fang, Zhongyuan Wang, and Ji-Rong Wen. An empirical study on eliciting and improving r1-like reasoning models. *arXiv preprint arXiv:2503.04548*, 2025b.
- Claude. Claude 3.7 sonnet. 2025. URL https://www.anthropic.com/news/claude-3-7-sonnet.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang,

Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948. 

- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models, 2023. URL https://arxiv.org/abs/2211.10435.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL https://arxiv.org/abs/2103.03874.
- Shijue Huang, Wanjun Zhong, Deng Cai, Fanqi Wan, Chengyi Wang, Mingxuan Wang, Mu Qiao, and Ruifeng Xu. Empowering self-learning of llms: Inner knowledge explicitation as a catalyst. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(23):24150–24158, Apr. 2025. doi: 10.1609/aaai. v39i23.34590. URL https://ojs.aaai.org/index.php/AAAI/article/view/34590.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report, 2024. URL https://arxiv.org/abs/2409.12186.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning, 2025. URL https://arxiv.org/abs/2503.09516.
- Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xiangru Peng, and Jiaya Jia. Step-dpo: Step-wise preference optimization for long-chain reasoning of llms, 2024. URL https://arxiv.org/abs/2406.18629.
- Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. Search-o1: Agentic search-enhanced large reasoning models, 2025a. URL https://arxiv.org/abs/2501.05366.
- Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl, 2025b. URL https://arxiv.org/abs/2503.23383.
- Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. Reft: Reasoning with reinforced fine-tuning, 2024. URL https://arxiv.org/abs/2401.08967.

518

519 520

521

522

523 524

525

526 527

528

529

530

531

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants, 2023. URL https://arxiv.org/abs/2311.12983.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025. URL https://arxiv.org/abs/2501.19393.

OpenAI. Learning to reason with llms, September 2024. URL https://openai.com/index/learning-to-reason-with-llms/.

OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. Openai o1 system card, 2024. URL https://arxiv.org/abs/2412.16720.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL https://arxiv.org/abs/2203.02155.

- Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. Logic-LM: Empowering large language models with symbolic solvers for faithful logical reasoning. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 3806–3824, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023. findings-emnlp.248. URL https://aclanthology.org/2023.findings-emnlp.248/.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL https://arxiv.org/abs/2412.15115.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023. URL https://arxiv.org/abs/2311.12022.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL https://arxiv.org/abs/2408.03314.
- Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. R1-searcher: Incentivizing the search capability in Ilms via reinforcement learning, 2025. URL https://arxiv.org/abs/2503.05592.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Kimi Team. Kimi k1.5: Scaling reinforcement learning with llms, 2025a. URL https://arxiv.org/abs/2501.12599.
- NovaSky Team. Sky-t1: Train your own o1 preview model within \$450. 2025b. URL https://novasky-ai.github.io/posts/sky-t1.
- OpenThoughts Team. Open Thoughts. https://open-thoughts.ai, January 2025c.
- Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025d. URL https://qwenlm.github.io/blog/qwq-32b/.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. Towards understanding chain-of-thought prompting: An empirical study of what matters. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2717–2739, Toronto, Canada,

- July 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.153. URL https://aclanthology.org/2023.acl-long.153/.
  - Hongru Wang, Yujia Qin, Yankai Lin, Jeff Z. Pan, and Kam-Fai Wong. Empowering large language models: Tool learning for real-world interaction. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, pp. 2983–2986, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400704314. doi: 10.1145/3626772.3661381. URL https://doi.org/10.1145/3626772.3661381.
  - Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song, Mingjie Zhan, and Hongsheng Li. Mathcoder: Seamless code integration in Ilms for enhanced mathematical reasoning, 2023b. URL https://arxiv.org/abs/2310.03731.
  - Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting, 2024. URL https://arxiv.org/abs/2402.10200.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 24824–24837. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper\_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.
  - Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Gang Fu, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Webdancer: Towards autonomous information seeking agency, 2025. URL https://arxiv.org/abs/2505.22648.
  - xAI. Grok. https://x.ai/, 2023. URL https://x.ai/. Large language model.
  - Yuxi Xie, Anirudh Goyal, Wenyue Zheng, Min-Yen Kan, Timothy P. Lillicrap, Kenji Kawaguchi, and Michael Shieh. Monte carlo tree search boosts reasoning via iterative preference learning, 2024. URL https://arxiv.org/abs/2405.00451.
  - An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024a.
  - An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement, 2024b. URL https://arxiv.org/abs/2409.12122.
  - Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: deliberate problem solving with large language models. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.

Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mammoth: Building math generalist models through hybrid instruction tuning, 2023. URL https://arxiv.org/abs/2309.05653.

Ge Zhang, Scott Qu, Jiaheng Liu, Chenchen Zhang, Chenghua Lin, Chou Leuang Yu, Danny Pan, Esther Cheng, Jie Liu, Qunshu Lin, Raven Yuan, Tuney Zheng, Wei Pang, Xinrun Du, Yiming Liang, Yinghao Ma, Yizhi Li, Ziyang Ma, Bill Lin, Emmanouil Benetos, Huan Yang, Junting Zhou, Kaijing Ma, Minghao Liu, Morry Niu, Noah Wang, Quehry Que, Ruibo Liu, Sine Liu, Shawn Guo, Soren Gao, Wangchunshu Zhou, Xinyue Zhang, Yizhi Zhou, Yubo Wang, Yuelin Bai, Yuhan Zhang, Yuxiang Zhang, Zenith Wang, Zhenzhu Yang, Zijian Zhao, Jiajun Zhang, Wanli Ouyang, Wenhao Huang, and Wenhu Chen. Map-neo: Highly capable and transparent bilingual large language model series, 2024. URL https://arxiv.org/abs/2405.19327.

Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang, Yifan Shao, Qichen Ye, Dading Chong, Zhiling Jin, Chenxuan Xie, Meng Cao, Yuxin Gu, Sixin Hong, Jing Ren, Jian Chen, Chao Liu, and Yining Hua. Browsecomp-zh: Benchmarking web browsing ability of large language models in chinese, 2025. URL https://arxiv.org/abs/2504.19314.

# A APPENDIX

# A.1 THE USE OF LARGE LANGUAGE MODELS (LLMS)

In preparing this paper, we employed large language models (LLMs) solely as writing assistants. Their role was limited to polishing grammar and improving fluency. They were not used for research ideation, methodological development, data analysis, or result interpretation. All scientific contributions, experimental designs, and findings presented in this work are entirely the result of the authors' own efforts.

#### A.2 RESULTS ON OTHER BASE MODEL

To further assess the effectiveness of ReTool, we conduct experiments using a code-oriented backbone, Qwen2.5-Coder-7B-Instruct (Hui et al., 2024). Following the same evaluation protocol as in the main results, we report the overall average accuracy on AIME2024, AIME2025, GSM8K, and MATH500. The results are presented in Table 2. We observe that applying our ReTool framework—through both supervised fine-tuning (SFT) and reinforcement learning (RL)—consistently yields substantial improvements across all benchmarks when built on top of Qwen2.5-Coder-7B-Instruct. These findings highlight the generalizability and effectiveness of ReTool across diverse base models.

Table 2: Results on different base model.

| Model                     | AIME2024 | AIME2025 | GSM8K | MATH500 |
|---------------------------|----------|----------|-------|---------|
| Qwen2.5-Coder-7B-Instruct | 10.0     | -        | 86.7  | 66.8    |
| + ReTool-SFT              | 14.69    | 17.29    | 90.67 | 77.2    |
| + ReTool-RL               | 46.04    | 32.60    | 92.08 | 84.55   |

#### A.3 THE PRACTICAL HELPFULNESS OF THE GENERATED CODE TOWARD THE FINAL SOLUTION.

To more rigorously evaluate the practical usefulness of generated code for solving the target tasks, we conducted a focused analysis of code snippets produced by our models before and after RL training. Specifically, we extracted all code snippets along with their corresponding contexts from the AIME2024 results and asked Claude4-Sonnet (Anthropic, 2025) to assess whether each snippet contributed meaningfully to solving the overall problem. The analysis reveals that, prior to RL training, 406 out of 471 code snippets (86.20%) were judged to be helpful, while 65 (13.80%) were not. After RL fine-tuning, the proportion of helpful code increased: 4,468 out of 5,010 snippets (89.18%) were rated as helpful, compared with 542 (10.82%) deemed unhelpful. These findings indicate not only a consistently high proportion of practically useful code but also a further improvement in code helpfulness following RL training.

# A.4 CI-POWERED REASONING VS. TEXT-BASED REASONING.

We present a case study to illustrate the distinction between CI-powered reasoning after reinforcement learning (RL) training and conventional text-based reasoning prior to RL training, as illustrated in Figure 7 in the Appendix. When faced with the same question, text-based reasoning relies on a "laborious" text-only calculation process, which is prone to numerical errors and often results in incorrect inference outcomes. In contrast, CI-powered reasoning substitutes this complex calculation process with concise code. This approach not only ensures computational accuracy through the assistance of an external code interpreter but also enables the model to focus more effectively on holistic reasoning strategies.

```
753
754
755
756
757
758
759
760
                 Question: From an unlimited supply of 1-cent coins, 10-cent coins, and 25-cent coins, ... Find the number of values of $N$
761
                 between $1$ and $1000$ inclusive for which the greedy algorithm succeeds.
762
                 Response: Alright, let's tackle this problem step by step using code to ensure accuracy. The goal is to determine for each value of
763
                 \( N \) from 1 to 1000 whether the greedy algorithm produces the minimum number of coins. If it does, we count that \( N \)...
764
                 <code>
                 ```python
765
                 def test_specific_N(N_values):
766
                   for N in N_values:
767
                     greedy_result = greedy(N)
768
                 test_results = test_specific_N(N_values_to_test)
769
                 print(test_results)
770
                 </code>
771
                 <interpreter>Traceback (most recent call last):
772
                 File \ ''/tmp/tmpbr3uj7zz/tmpdmj39k6o.py'', line \ 11, in < module > test\_results = test\_specific\_N(N\_values\_to\_test)
773
                 File "/tmp/tmpbr3uj7zz/tmpdmj39k6o.py", line 5, in test_specific_N greedy_result = greedy(N)
                 NameError: name 'greedy' is not defined </interpreter>
774
775
                 Oops, the functions need to be defined in the same scope. Let's correct that:
776
                 <code>
777
                 ```python
778
                 def greedy(N):
                   denominations = [25, 10, 1]
779
780
                 def test_specific_N(N_values):
781
                   for N in N values:
                     greedy_result = greedy(N)
782
783
                 test_results = test_specific_N(N_values_to_test)
784
                 print(test_results)
785
                 </code>
786
                 787
                 <answer>\boxed{610} </answer>
788
789
```

Figure 6: The case of "aha moment" about code self-correction.

836

```
806
807
                Question: Let $N$ be the greatest four-digit positive integer with the property that whenever one of its digits is changed to $1$, the resulting number is divisible
808
                by $7$. Let $Q$ and $R$ be the quotient and remainder, respectively, when $N$ is divided by $1000$. Find $Q+R$.
809
                Text-based Reasoning Before RL: Okay, so I need to find the greatest four-
                                                                                                      CI-powered Reasoning After RL: Alright, let's tackle this problem step by
                digit number N such that if any one of its digits is replaced by 1, the resulting
                                                                                                      step. The goal is to find the greatest four-digit positive integer \backslash (\ N\ \backslash) such that
810
                number is divisible by 7. Then, I have to find Q and R when N is divided by
                                                                                                      changing any one of its digits to 1 results in a number divisible by 7. Then, we
811
                1000, and then compute Q + R. Alright, let's break this down step by step.
                                                                                                      need to find \( Q + R \) where \( Q \) and \( R \) are the quotient and remainder
                Therefore, each of these equations can be rearranged to:
                                                                                                      when \( N \) is divided by 1000.
812
                                                                                                      First, let's compute the values of \( 1000 \mod 7 \), \( 100 \mod 7 \), \( 10 \mod 7 \)
                1. N \equiv -(1 - A)*1000 \mod 7
                2. N \equiv -(1 - B)*100 \mod 7
813
                                                                                                      7 \), and \( 1 \mod 7 \):
                3. N \equiv -(1 - C)*10 \mod 7
                                                                                                      <code>
814
                4. N \equiv -(1 - D)*1 \mod 7
                                                                                                         python
815
                Therefore, all four expressions must be equal to N mod7. Therefore, we can set
                                                                                                      # Compute modular values
                them equal to each other: -(1 - A)*1000 \equiv -(1 - B)*100 \equiv -(1 - C)*10 \equiv -(1 - C)
                                                                                                      mod_1000 = 1000 % 7
816
                                                                                                      mod_100 = 100 % 7
                Wait, let's check B, C, D for A=9:
                                                                                                      mod 10 = 10 % 7
817
                B=3A -2 mod7=3*9 -2=27 -2=25 mod7=25 -7*3=25-21=4
                                                                                                      mod_1 = 1 % 7
818
                C=2A -1 mod7=2*9 -1=18 -1=17 mod7=17 -7*2=17-14=3
                                                                                                      print(f"1 mod 7 = {mod 1}")
                D=6A -5 mod7=6*9 -5=54 -5=49 mod7=49 -7*7=49-49=0
819
                So, B=4, C=3, D=0. Therefore, N=9430...
820
                Let me compute each of these modulo7:
                                                                                                      </code>
                1. 1000 ≡6 mod7, so 6 + 100*B + 10*C + D ≡0 mod7
                                                                                                      <interpreter>1000 mod 7 = 6 100 mod 7 = 2 10 mod 7 = 3 1 mod 7 =
821
                But 100*B ≡2*B mod7,
                                                                                                      1</interpreter>
                10*C \equiv 3*C \mod 7, so:
822
                6 + 2B + 3C + D \equiv 0 \mod 7
                                                                                                      We need to solve these equations for digits \langle (a, b, c, d \rangle) where \langle (a \rangle) is from 1
823
                2. 1000*A \equiv 6A \mod 7, 100 \equiv 2 \mod 7,
                                                                                                      to 9 and \setminus( b, c, d \setminus) are from 0 to 9. Let's use code to iterate through possible
                                                                                                      values and find the valid \setminus ( N \setminus).
824
                6A + 2 + 3C + D \equiv 0 \mod 7
                                                                                                       <code>
825
                3. 1000*A + 2B + 3 + D \equiv 0 \mod 7 \text{ (since } 100*B \equiv 2B, 10*1 \equiv 3, 10*C \equiv 3C)
                                                                                                      ```python
def find_valid_n():
                Wait, no. Let me recompute each equation properly.
826
  for a in range(9, 0, -1):
                Equation1: 1BCD ≡0 mod7
827
                Which is 1000*1 + 100*B + 10*C + D \equiv 0 \mod 7 \ 1000 \equiv 6, \ 100 \equiv 2, \ 10 \equiv 3, \dots
                Let me tabulate A from 1 to 9:
   return N
828
                A | B=5-4A mod7 | C=6-5A mod7 | D=2-A mod7 | Valid?
829
                 1 | 5-4=1 mod7=1 |6-5=1 mod7=1 |2-1=1 |Yes...
  N = find_valid_n()
                **Result**
  Q = N // 1000
830
                - \( N = 5624 \)
  R = N % 1000
                - Dividing by 1000 gives \( Q = 5 \) and \( R = 624 \)
831
  result = O + R
  print(f"N = {N}, Q = {Q}, R = {R}, Q + R = {result}")
                - \ (Q + R = 5 + 624 = 629 \)
832
                \[ \boxed{629} \]
833
  <interpreter>N = 5694, Q = 5, R = 694, Q + R = 699 /interpreter> ...
834
  Thus, the final answer is confirmed to be 699.
   <answer>\boxed{699} </answer>
835
```

Figure 7: Case of CI-powered reasoning vs. Text-based reasoning.

Template prompt for Reinforcement Learning Rollout Solve the following problem step by step. You now have the ability to selectively write executable Python code to enhance your reasoning process. The Python code will be executed by an external sandbox, and the output (wrapped in <interpreter>output</interpreter>) can be returned to aid your reasoning and help you arrive at the final answer. The Python code should be complete scripts, including necessary imports. **Code Format:** Each code snippet is wrapped with <code> ```python code snippet </code> **Answer Format:** The last part of your response should be: <answer>\boxed{'The final answer goes here.'}</answer> **User Question:** {question} Assistant: Figure 8: Template prompt for ReTool rollout. 

your calculation while still maintaining the reasoning process.

897 898

# 899 900

901 902

903 904

905 906

907

908 909

910 911

912 913

914 915

916

917 918 919

920 921

922 923

924 925

926 927

928 929 930

931 932 933

You are a helpful AI assistant. Initially, when solving a question, you would need to think step by step, without the ability to use code for calculation. Now, you have the capability to write code to use the code interpreter for calculation. The code will be executed by a sandbox, and the result can be returned to enhance your reasoning process. You can now leverage code to enhance

The thinking process can have multiple code snippets. Each code snippet is wrapped with: <code>

```python

code snippet

</code>, and should be executable. The returned result is wrapped with <interpreter> execution results \texttt{</interpreter>}.

**Template Prompt for Data Curation** 

Modify the original thinking process to make it more accurate by replacing manual calculation steps that can benefit from code execution with the corresponding code snippets and their interpreter's execution results. The core reasoning logic from the original thinking process, including any unsuccessful attempts, should remain unchanged. You should only replace the necessary manual calculation steps with code and interpreter's execution results, without altering the rest tokens of the thinking process. Wrap the revised thinking process within <revised\_thinking\_process> and </revised\_thinking\_process>}.

#### **User Question:**

{question}

## Original Thinking Process (without code interpreter's support):

<original\_thinking\_process> {original\_response} </original\_thinking\_process>

#### Details:

- 1. Identify sections where code execution could speed up the reasoning process or make the calculation more accurate.
- 2. Replace the manual calculation steps with code snippets and the corresponding interpreter's execution results.
- 3. Keep the logical flow of the reasoning process intact, including any failed exploration attempts that were part of the initial
- 4. The code snippets should be complete scripts, including necessary imports, and should not contain markdown symbols like <code>

```python

code snippet

</code>.

- 5. Outputs in the code snippets must explicitly call the **print** function.
- 6. Execution results should match the model's output exactly, with no extra or missing tokens.
- 7. If the Original Thinking Process does not include an <answer> section at the end, please add it in the Revised Thinking Process: <answer> \boxed{'The final answer goes here.'} </answer>

# Revised Thinking Process (With code interpreter's support):

Figure 9: Template prompt for data curation.