

MENTOR: A Reinforcement Learning Framework for Distilling Tool Use in Small Models via Teacher-Optimized Rewards

Anonymous ACL submission

Abstract

Distilling the tool-using capabilities of large language models (LLMs) into smaller, more efficient small language models (SLMs) is a key challenge for their practical application. The predominant approach, supervised fine-tuning (SFT), suffers from poor generalization as it trains models to imitate a static set of teacher trajectories rather than learn a robust methodology. While reinforcement learning (RL) offers an alternative, the standard RL methods using simple outcome-based reward fail to effectively guide SLMs, causing them to struggle with inefficient exploration and adopt suboptimal strategies. To address these distinct challenges, we propose MENTOR, a framework that synergistically combines RL with teacher-guided distillation. Instead of simple imitation, MENTOR employs an RL-based process to learn a more generalizable policy through exploration. In addition, to address insufficient process guidance, it uses a teacher’s reference trajectory to construct a composite teacher-guided reward that provides fine-grained guidance. Extensive experiments demonstrate that MENTOR significantly improves the cross-domain generalization and strategic competence of SLMs compared to both SFT and standard RL baselines.

1 Introduction

The augmentation of large language models (LLMs) with external tools, such as code interpreters and retrieval APIs, has enabled them as advanced agents capable of handling complex reasoning tasks (Yao et al., 2023; Paranjape et al., 2023; Wang et al., 2024b; Singh et al., 2025). However, the high inference costs of these large-scale models hinder their practical use. This challenge has motivated a line of research focused on smaller language models (SLMs), with the goal of preserving the tool-assisted problem-solving capabilities of larger models (Gao et al., 2023; Gou et al., 2024; Qiu et al., 2025).

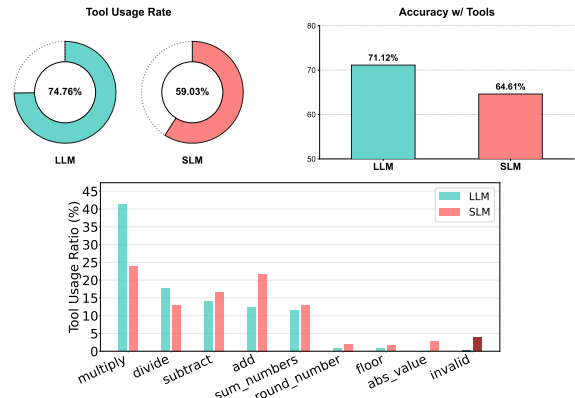


Figure 1: Comparative analysis of tool usage rates and their corresponding performance, as well as tool invocation patterns between LLMs and SLMs.

Knowledge distillation has been steadily used for transferring such advanced abilities from a larger **teacher** model to a smaller **student** model (Ho et al., 2023; Chenglin et al., 2024; Dai et al., 2024b, 2025; Song et al., 2025; Liao et al., 2025). These studies primarily employ supervised fine-tuning (SFT), where a student model is trained to replicate teacher-generated problem-solving trajectories (Liu et al., 2024; Kang et al., 2025; Lyu et al., 2025). However, the reliance of these methods on a static dataset presents a fundamental scalability issue, as it is impossible to generate trajectories for every scenario the model will face (Luo et al., 2025b; Sun et al., 2025; Yin et al., 2025). As a result, these models often fail to generalize to new tasks, necessitating more adaptive frameworks for teaching tool use.

As an alternative to this imitation-based approach, some studies have explored reinforcement learning (RL) through iterative self-refinement (Trung et al., 2024; Wu et al., 2025b). To guide the models in learning from their solutions, prior works use simple reward signals, such as correctness of the final answer or tool invocation (Yu et al., 2024; Singh et al., 2025; Wu et al., 2025a). However, relying on this naive outcome sig-

069 nal is insufficient for SLMs, as their inefficient use
070 of tools often leads them to adopt suboptimal strate-
071 gies. As shown in Figure 1, which compares the
072 tool-use strategies of a teacher LLM and a student
073 SLM, SLMs are less effective at leveraging tools:
074 the Teacher utilized tools in 74.8% of samples
075 (71.12% accuracy), whereas the Student utilized
076 tools in only 59.0% of samples (64.61% accuracy).
077 This is linked to suboptimal strategies in SLMs,
078 as also evidenced by their distinct tool-invocation
079 patterns and tendency to issue invalid tool calls.
080 This highlights the importance of fine-grained feed-
081 back to guide the agent toward an effective tool-use
082 strategy.

083 To address these distinct challenges of both SFT-
084 based distillation and standard RL, we propose
085 a framework that synergistically combines rein-
086 forcement learning with teacher-guided distillation.
087 We introduce MENTOR (Model ENhancement via
088 Teacher-Optimized Rewards), a novel approach de-
089 signed to leverage the strengths of each approach.
090 By employing an **RL-based distillation** process,
091 our framework moves beyond the simple imita-
092 tion of SFT, allowing the student model to learn
093 a more generalizable policy through exploration.
094 Simultaneously, to address the challenge of insuf-
095 ficient process guidance that hinders SLMs, we
096 design a **teacher-guided reward**, which utilizes
097 the teacher’s trajectory as a reference to construct
098 a composite reward signal. This reference-based re-
099 ward provides the fine-grained guidance necessary
100 to steer the SLM towards efficient tool-use strate-
101 gies and prevent it from converging to suboptimal
102 policies. Our code is publicly available¹.

103 To summarize, our key contributions are:

- 104 • To address the limitations of SFT-based approach,
105 we propose MENTOR. Unlike SFT, which of-
106 ten leads to superficial imitation, our RL-based
107 framework uses exploration to internalize the
108 teacher’s tool-use strategy, ensuring the student
109 learns a generalizable methodology.
- 110 • To prevent SLMs from adopting suboptimal
111 strategies due to non-comprehensive signals, we
112 propose a teacher-optimized reward that lever-
113 ages the teacher’s trajectory for fine-grained su-
114 pervision, ensuring the student internalizes the
115 correct methodology.
- 116 • We demonstrate through extensive experiments

¹<https://anonymous.4open.science/r/MENTOR-F6E7/>

117 that MENTOR significantly improves the cross-
118 domain generalization and strategic competence
119 of SLMs compared to both SFT and standard
120 outcome-based RL baselines.

121 2 Related Work

122 2.1 Tool-Augmented Language Models

123 A significant line of research enhances the reason-
124 ing of LLMs by augmenting them with external
125 tools, such as code interpreters and retrieval APIs.
126 The approaches for integrating these tools are di-
127 verse, ranging from prompting-based code genera-
128 tion (Gao et al., 2023; Paranjape et al., 2023; Inaba
129 et al., 2023; Huang et al., 2024) to explicitly train-
130 ing models to invoke tools as part of their reasoning
131 process (Schick et al., 2023; Kong et al., 2023; Qian
132 et al., 2024). Recognizing that effective tool use
133 is an inherently strategic process, recent work has
134 focused on learning optimal tool-invocation pat-
135 terns through RL (Yu et al., 2024; Feng et al., 2025;
136 Qian et al., 2025). The core challenge in training
137 tool-augmented models is maintaining a robust and
138 generalizable tool-use policy across long reasoning
139 trajectories.

140 2.2 Reasoning Distillation through SFT

141 Previous studies have mainly trained student mod-
142 els to clone teacher-generated trajectories as a
143 method for transferring tool-use capabilities to
144 SLMs (Liu et al., 2024; Kang et al., 2025; Lyu
145 et al., 2025). However, SFT-based distillation faces
146 a critical scalability challenge (Sun et al., 2025;
147 Yin et al., 2025), as it is infeasible to curate a static
148 dataset that covers every possible scenario. A key
149 limitation that arises from this data dependency is
150 that models learn to mimic the superficial format
151 of a reasoning trajectory to produce a correct final
152 answer, without internalizing the underlying logi-
153 cal process (Kandpal et al., 2023; Dai et al., 2024a;
154 Li et al., 2025). By contrast, MENTOR employs
155 an **RL-based distillation** framework, where the
156 teacher’s trajectory is used not for direct imitation,
157 but serves as a reference to guide an exploratory
158 learning process aimed at developing a more **gen-
159 eralizable** problem-solving method.

160 2.3 Reinforcement Learning and Reward 161 Design

162 Reinforcement learning (RL) (Kaelbling et al.,
163 1996) offers a powerful alternative to SFT’s
164 imitation-based learning, as its exploratory nature

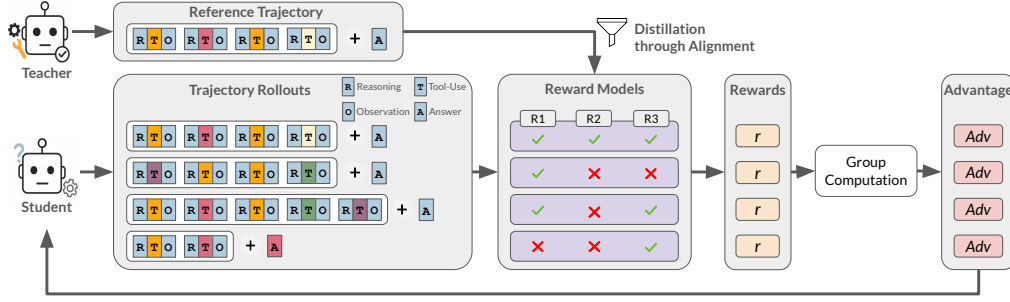


Figure 2: Overview of the MENTOR training framework. A problem-solving trajectory (τ) consists of a sequence of Reasoning (R), Tool-Use (T), and Observation (O), and a final Answer (A). Each student rollout is evaluated by a set of reward models, which generate a reward signal by aligning the student’s actions against the teacher’s reference trajectory.

can, in principle, discover more generalizable policies. Foundational algorithms like Proximal Policy Optimization (PPO) established the paradigm of fine-tuning models through policy optimization (Schulman et al., 2017). More recent advancements, such as Group Relative Policy Optimization (GRPO), have adapted this approach for complex reasoning tasks, demonstrating its effectiveness in fostering robust, self-corrective behaviors (Shao et al., 2024).

Many of these RL approaches have been demonstrated on highly capable LLMs for reasoning with tool use. Due to their strong intrinsic abilities, these models can often discover effective policies even when guided by simple outcome rewards, such as a binary signal for final answer correctness or successful tool invocation (Yu et al., 2024; Feng et al., 2025; Singh et al., 2025; Qian et al., 2025; Wu et al., 2025a). However, SLMs are significantly less efficient explorers (Wei et al., 2022; Xiong et al., 2024). With only outcome-based rewards, they struggle to link their actions to the final outcome and tend to converge to suboptimal policies. To address this, our work leverages the teacher’s trajectory to construct a **teacher-guided reward** signal, providing the fine-grained guidance necessary to steer the SLM’s exploration.

3 MENTOR: Model Enhancement via Teacher-Optimized Reward

In this section, we introduce MENTOR, a framework that leverages reinforcement learning (RL) to distill a tool-calling policy from a large teacher model to a smaller student model. The overall process is illustrated in Figure 2. For each input, the teacher model first generates a reference reasoning trajectory that exemplifies a successful problem-solving process. Concurrently, the student model generates multiple exploratory rollouts to sample

different reasoning paths. We then apply the Group Relative Policy Optimization (GRPO) algorithm to refine the student’s policy. By leveraging a teacher-guided reward signal, MENTOR ensures the student internalizes strategic principles rather than merely memorizing steps. We describe this process in the following subsections.

3.1 Reference Trajectory Generation

First, we use a large teacher model (π_{teacher}) to generate reference reasoning trajectories (τ). Each trajectory consists of a sequence of reasoning (r), tool-use (t), and observation (o) steps. We use a training dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, where x_i represents a question and y_i is its corresponding ground-truth answer. For each question x , the teacher model generates a final answer \hat{y} and a reasoning trajectory τ using an instruction prompt (I) as follows:

$$O^{(t)} = (\tau^{(t)}, \hat{y}^{(t)}) \sim \pi_{\text{teacher}}(\cdot | x, I), \text{ where} \quad (1)$$

$$\tau^{(t)} = \langle (r_1, t_1, o_1), \dots, (r_{L_\tau}, t_{L_\tau}, o_{L_\tau}) \rangle. \quad (2)$$

3.2 Adapting GRPO for Reasoning Distillation

We adapt the GRPO algorithm to distill the teacher’s problem-solving methodology by configuring the teacher’s successful trajectory as a high-reward target. This reward-driven setup trains the student to internalize the teacher’s strategic tool-use policy, rather than merely imitating a fixed sequence of actions. This process is detailed in Algorithm 1.

Generating Rollouts For each question x_i in our training set, we perform two simultaneous generation steps. First, we retrieve the corresponding reference trajectory, $(\tau^{(t)}, \hat{y}^{(t)})$, which was generated by the teacher model as described in subsection 3.1. Then, we use the current student model (π_{old}) to

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{x_i \sim \mathcal{D}, \{O_j^{(s)}\}_{j=1}^G \sim \pi_{old}(\cdot|x_i)} \left[\frac{1}{G} \sum_{j=1}^G \left(\sum_{k=1}^{|\tau_j^{(s)}|} \mathbb{I}(\tau_{j,k}) \right)^{-1} \right. \\ \left. \sum_{k=1}^{|\tau_j^{(s)}|} \min \left(\frac{\pi_{\theta}(\tau_{j,k}|\tau_{j,<k},x_i)}{\pi_{old}(\tau_{j,k}|\tau_{j,<k},x_i)} \hat{A}_{j,k}, \text{clip} \left(\frac{\pi_{\theta}(\tau_{j,k}|\tau_{j,<k},x_i)}{\pi_{old}(\tau_{j,k}|\tau_{j,<k},x_i)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{j,k} \right) \cdot \mathbb{I}(\tau_{j,k}) - \beta \mathbb{D}_{KL}[\pi_{\theta}||\pi_{ref}] \right] \quad (3)$$

Algorithm 1 Training with GRPO

Require: Student model π_{θ} , old student model π_{old} , Teacher model $\pi_{teacher}$, task dataset \mathcal{D} , group size G , indicator function \mathbb{I}

- 1: **for** each training iteration **do**
- 2: **for** each question x_i **do**
- 3: Generate reference trajectory $O^{(t)}$ from $\pi_{teacher}$
- 4: Sample G rollouts $\{O_1^{(s)}, \dots, O_G^{(s)}\}$ from π_{old}
- 5: **for** each rollout $O_j^{(s)}$ **do**
- 6: Compute outcome rewards $R(O_j^{(s)}, O^{(t)})$
- 7: **end for**
- 8: Compute groupwise advantages $\hat{A}_{j,k}$ for all $O_j^{(s)}$
- 9: Apply indicator \mathbb{I} to mask tool output tokens
- 10: Compute GRPO loss \mathcal{L}_{GRPO} and update $\pi_{student}$
- 11: **end for**
- 12: **end for**

generate a group of G rollouts. This step is crucial as it allows the student to explore diverse reasoning paths and tool-use strategies for the same problem, providing the varied data needed for the GRPO comparison.

Student Policy Optimization The core idea of our proposed method is to update a policy by aligning a high-quality reference against a group of sampled candidates. We optimize the tool-use policy of the student model based on the reference trajectory generated by the teacher. We optimize the student model by maximizing the objective function as shown in Equation 3, where ϵ and β are hyperparameters for clipping and KL regularization, respectively. The advantage, $\hat{A}_{j,k}$, is computed from the relative rewards within the sample group. The reference policy π_{ref} used for KL regularization is the initial student model before RL training. $\mathbb{I}(\tau_{j,k})$ is an indicator function used for loss masking, which equals 1 if $\tau_{j,k}$ is an LLM-generated token, and 0 otherwise.

3.3 Teacher-Guided Reward Design

The central challenge of our work lies in the reward design: how to distill a generalizable problem-solving methodology, rather than merely rewarding correct answers. An outcome-based reward is insufficient because it lacks the granular supervision needed to guide the student through the teacher’s strategic problem-solving process. To address this, we provide a more fine-grained signal that also evaluates the reasoning process itself by designing a composite reward mechanism with three key

components. The total reward for a given student output, $R(O^{(s)}, O^{(t)})$, is defined as:

$$R(O^{(s)}, O^{(t)}) = w_c R_c + w_a R_a + w_v R_v, \quad (4)$$

where w_c , w_a , and w_v are hyperparameters that balance the contribution of each reward component.

Correctness Reward This component evaluates whether the final answer derived by the student model ($\hat{y}^{(s)}$) is consistent with the result produced by the teacher model ($\hat{y}^{(t)}$). This provides a direct signal indicating if the student’s overall reasoning process concludes with the same outcome as the teacher’s successful demonstration. For our training, we only use reference trajectories where the teacher’s answer matches the ground truth (y). The reward is defined as:

$$R_c = \begin{cases} 1 & \text{if } \hat{y}^{(s)} = \hat{y}^{(t)} \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

Teacher-Alignment Reward To provide an intermediate signal that guides the student towards the teacher’s problem-solving strategy, we introduce this reward. This component encourages the student to select the same set of tools as the teacher, which is a crucial aspect of learning the overall methodology. The reward is assigned only if the set of tool calls made in the student’s trajectory, $\tau^{(s)}$, is identical to the set of tool calls in the teacher’s trajectory, $\tau^{(t)}$. This is defined as:

$$R_a = \begin{cases} 1 & \text{if } \{t_1^{(s)}, \dots, t_L^{(s)}\} = \{t_1^{(t)}, \dots, t_L^{(t)}\} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

, where each $t_i^{(s)}$ in $\{t_1^{(s)}, \dots, t_L^{(s)}\}$ are from $\tau^{(s)}$, and each $t_i^{(t)}$ in $\{t_1^{(t)}, \dots, t_L^{(t)}\}$ are from $\tau^{(t)}$.

Tool Validation Reward A primary inefficiency we observed (in Figure 1) was the student model’s tendency to generate invalid tool calls. These actions, which result in execution errors from the tool interpreter, immediately derail the reasoning process and are a significant source of suboptimal performance for the SLM. To directly penalize this behavior and guide the student towards the valid, error-free trajectories demonstrated by the teacher,

we introduce the tool validation reward. This binary reward is 1 only if every tool call within the student’s trajectory, $\tau_{student}$, executes successfully without raising an error:

$$R_V = \begin{cases} 1 & \forall o_i^{(s)} \in \tau^{(s)} \text{ is valid} \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

4 Experiments and Results

4.1 Experimental Setup

Training Domain Selection We strategically select mathematical reasoning as the primary training domain to enable the model to explore various problem-solving approaches. Its inherent difficulty gradient drives robust exploration beyond simple policies, mitigating the risk of suboptimal convergence (Zhou et al., 2023; Wang et al., 2024a; Chen et al., 2025b; Luo et al., 2025a). This choice is also supported by prior work showing that mathematical ability is a transferable skill that provides a strong foundation for generalizable problem-solving (Wang et al., 2025; Huang et al., 2025). We use the AceReason-Math dataset for training (Chen et al., 2025b). Details of our training dataset are provided in Appendix A.1.

Models and Tools Our teacher model is Qwen3-235B-Thinking, chosen for its strong tool-use capabilities. We use four student models to evaluate our method across different scales: Qwen3 (8B and 1.7B) and Qwen2.5 (7B and 1.5B). To augment the models with tool-use capability, we implement a sandbox for running Python code on a remote server. Our implementation code is publicly available. The details of the model versions and tools are in Appendix A.2 and Appendix C, respectively.

Baselines To evaluate our proposed framework, we compare MENTOR against three distinct baselines that represent different training approaches. **1) Vanilla SLM:** The base instruction-tuned model without any further training. This serves as a lower bound for performance. **2) SFT:** This baseline represents the predominant distillation method, where the SLM is fine-tuned on expert-generated trajectories using supervised fine-tuning. **3) Standard RL:** To isolate the benefit of our teacher-guided reward design, we include a Standard RL method as a baseline. This agent is trained using the same RL framework as our method, but with a simple reward signal (R_c) based only on the final outcome.

Task Type	Dataset Name	Description	Size
Math Reasoning	Math-Forge-Hard	College	500
	Omni-MATH-512 (Gao et al., 2024)	Olympiad	512
	AIME24 (AI-MO, 2024)	Olympiad	30
	AIME25	Olympiad	30
	amc23 (AMC, 2023)	Olympiad	40
	minervamath (Lewkowycz et al., 2022)	College	272
Tool-Calling	BFCL v4 (Patil et al., 2025)	Tool-call	5088
Factual Reasoning	HotPotQA (Yang et al., 2018)	2-hop QA	2000
	2WikiMultiHopQA (Ho et al., 2020)	2-hop QA	2000
	Bamboogle (Press et al., 2023)	2-hop QA	125

Table 1: Benchmarks categorized by in-domain (*Mathematical Reasoning*) and out-of-domain (*Tool-Calling* and *Factual Reasoning*) tasks and their test data size.

Benchmarks We evaluate our framework on a set of in-domain and out-of-domain tasks, detailed in Table 1, to measure both task-specific performance and generalization ability. For **In-Domain Tasks**, our agent is trained and evaluated on a collection of mathematical reasoning benchmarks that provide a natural gradient of difficulty. This includes the widely-used MATH dataset (Hendrycks et al., 2021), the tool-centric Omni-MATH-512 (Gao et al., 2024), and several Olympiad-level datasets such as AIME24, AIME25, amc23, and minervamath. For **Out-of-Domain Tasks**, we evaluate the trained agent on tasks requiring tools unseen during training to assess zero-shot generalization. To test *retrieval-based QA*, we reframe multi-hop QA datasets (HotpotQA (Yang et al., 2018), 2WikiMultiHopQA (Ho et al., 2020), Bamboogle (Press et al., 2023)) as a tool-use problem where the agent is provided with a search(query) tool and must learn to call it effectively. To test broader capabilities, we also use the general *Tool-Calling* benchmark BFCL v4 (Patil et al., 2025), which involves a diverse set of novel tools. We provide details in Appendix A.3.

Evaluation Metrics We evaluate the performance on all tasks using Exact Match (EM). To measure the accuracy (Acc), we consider a prediction to be correct only if it exactly matches one of the ground-truth answers after normalization. We also quantify policy alignment using an **alignment score (AS)** based on the Jensen-Shannon divergence. This score measures the divergence between a model’s tool usage distribution and the teacher’s reference distribution from Figure 1. Further details are in Appendix A.4.

Implementation Details The reinforcement learning framework is built on verl (Sheng et al., 2024). The agent is trained for two epochs on the combined training splits of our training dataset. We employ LoRA (Hu et al., 2021) for supervised fine-

Model	Method	Math (Acc)							Overall
		Math-Forge	Omni-MATH	aime24	aime25	amc23	minervamath		
Qwen 3	235B-Vanilla	68.00	26.00	46.67	40.00	85.00	46.69	52.06	
Qwen 2.5	1.5B-Vanilla	5.18 (± 0.11)	1.23 (± 0.18)	0.00 (± 0.00)	0.00 (± 0.00)	2.50 (± 2.04)	1.29 (± 0.56)	1.70	
	1.5B-SFT	5.79 (± 0.09)	1.62 (± 0.23)	1.00 (± 2.25)	1.33 (± 1.72)	5.75 (± 2.90)	3.53 (± 0.43)	3.17	
	1.5B-Standard RL	18.38 (± 0.09)	3.46 (± 0.18)	2.67 (± 3.06)	3.00 (± 2.92)	8.50* (± 2.93)	4.12 (± 0.34)	6.69	
	1.5B-MENTOR	18.84* (± 0.13)	5.64* (± 0.28)	10.00* (± 3.52)	6.67* (± 2.72)	9.75* (± 2.49)	8.42* (± 0.44)	9.89	
	7B-Vanilla	36.87 (± 1.04)	8.60 (± 1.28)	6.66 (± 2.22)	5.33* (± 2.81)	44.00 (± 3.16)	26.48 (± 0.96)	21.32	
Qwen 3	7B-SFT	37.64 (± 1.25)	11.14 (± 1.99)	9.40* (± 2.65)	7.33* (± 3.44)	43.00 (± 4.53)	26.38 (± 2.64)	22.48	
	7B-Standard RL	50.25 (± 1.32)	12.94 (± 1.21)	12.00* (± 3.58)	6.66* (± 2.22)	45.50 (± 3.07)	29.40 (± 1.70)	26.13	
	7B-MENTOR	55.42* (± 1.96)	14.72* (± 2.21)	12.33* (± 3.52)	7.33* (± 2.63)	50.10* (± 2.70)	31.86* (± 1.71)	28.63	
	1.7B-Vanilla	58.70 (± 0.11)	16.28 (± 0.21)	25.36 (± 2.32)	13.01 (± 2.44)	58.00 (± 4.22)	28.85 (± 0.47)	33.37	
	1.7B-SFT	59.99 (± 0.12)	16.65 (± 0.21)	31.00 (± 2.25)	16.02 (± 2.10)	63.75 (± 3.17)	30.60 (± 0.52)	36.34	
Qwen 3	1.7B-Standard RL	60.18 (± 0.07)	17.78 (± 0.28)	34.34 (± 3.13)	20.00 (± 2.22)	67.75 (± 2.75)	31.67 (± 0.48)	38.62	
	1.7B-MENTOR	60.66* (± 0.11)	20.74* (± 0.22)	36.67* (± 2.50)	24.00* (± 4.08)	71.00* (± 3.94)	32.30* (± 0.38)	40.90	
	8B-Vanilla	65.00 (± 0.12)	19.46 (± 0.24)	29.32 (± 2.64)	20.65 (± 3.44)	57.75 (± 2.49)	42.55 (± 0.58)	39.12	
	8B-SFT	65.35 (± 0.07)	20.45 (± 0.21)	31.65 (± 3.23)	25.65 (± 4.17)	62.50 (± 2.89)	43.22 (± 0.31)	41.47	
	8B-Standard RL	65.48 (± 0.04)	22.03 (± 0.26)	35.33 (± 3.22)	30.33 (± 3.99)	75.50 (± 2.58)	43.49 (± 0.25)	45.19	
8B-MENTOR	66.50* (± 0.11)	24.10* (± 0.25)	39.00* (± 3.53)	37.00* (± 3.31)	79.00* (± 2.69)	43.94* (± 0.44)	48.26		

Model	Method	BFCL-v4 (Acc)					RAG (EM)			
		Non-Live	Multi-Turn	Live	Agentic	Overall	Bamboogle	2WikiMultiHopQA	HotpotQA	Overall
Qwen 3	235B-Vanilla	87.62	51.88	82.68	18.83	50.13	41.60	42.50	34.60	39.57
Qwen 2.5	1.5B-Vanilla	68.88 (± 1.15)	1.09 (± 0.34)	58.99 (± 0.54)	2.21 (± 0.44)	24.00	0.30 (± 0.31)	3.20 (± 0.55)	1.47 (± 0.60)	1.66
	1.5B-SFT	69.96 (± 1.20)	1.58 (± 0.42)	60.72 (± 0.36)	2.65 (± 0.29)	24.59	0.85 (± 0.53)	3.91 (± 0.47)	2.52 (± 0.57)	2.43
	1.5B-Standard RL	71.38 (± 0.63)	1.98 (± 0.50)	61.46 (± 0.45)	3.09 (± 0.39)	25.12	6.80 (± 0.49)	10.19 (± 0.41)	5.79 (± 0.66)	7.59
	1.5B-MENTOR	72.66* (± 0.62)	2.35* (± 0.41)	62.03* (± 0.39)	3.83* (± 0.40)	25.70	7.36* (± 0.55)	11.31* (± 0.54)	6.50* (± 0.40)	8.39
	7B-Vanilla	72.04 (± 0.33)	3.90 (± 1.43)	63.78 (± 0.91)	4.09 (± 1.86)	26.38	14.88 (± 2.48)	14.20 (± 0.58)	13.08 (± 1.00)	14.05
Qwen 3	7B-SFT	72.05 (± 0.32)	4.02 (± 1.08)	63.79 (± 0.79)	4.27 (± 1.76)	26.50	16.48 (± 1.74)	13.87 (± 0.53)	13.36 (± 0.70)	14.57
	7B-Standard RL	82.13 (± 0.09)	14.14* (± 1.69)	72.22 (± 0.30)	8.15 (± 1.15)	32.94	22.04 (± 1.06)	17.57 (± 0.55)	17.07 (± 0.65)	18.89
	7B-MENTOR	82.35* (± 0.21)	14.27* (± 1.74)	72.81* (± 0.43)	9.30* (± 1.34)	33.52	23.56* (± 1.24)	19.39* (± 0.40)	19.24* (± 0.56)	20.73
	1.7B-Vanilla	80.89 (± 0.70)	10.47 (± 0.52)	69.90 (± 0.22)	3.98 (± 0.30)	29.81	14.61 (± 0.53)	18.87 (± 0.31)	16.11 (± 0.36)	16.53
	1.7B-SFT	81.55 (± 0.50)	10.45 (± 0.46)	70.36 (± 0.41)	4.51 (± 0.39)	30.13	15.95 (± 0.54)	19.69 (± 0.51)	17.00 (± 0.43)	17.45
Qwen 3	1.7B-Standard RL	83.00 (± 0.31)	11.50 (± 0.38)	70.77 (± 0.31)	4.73 (± 0.46)	30.65	16.83 (± 0.55)	20.36 (± 0.53)	17.35 (± 0.29)	18.18
	1.7B-MENTOR	82.76* (± 0.35)	12.02* (± 0.46)	71.05* (± 0.34)	5.52* (± 0.60)	31.20	18.28* (± 0.55)	21.37* (± 0.60)	18.03* (± 0.48)	19.23
	8B-Vanilla	87.28 (± 0.49)	35.51 (± 0.47)	80.16 (± 0.29)	9.89 (± 0.33)	41.35	32.13 (± 0.64)	35.88 (± 0.61)	27.47 (± 0.61)	31.83
	8B-SFT	87.91 (± 0.48)	36.92 (± 0.60)	80.59 (± 0.33)	10.58 (± 0.44)	42.16	34.65 (± 0.31)	37.18 (± 0.36)	30.59 (± 0.44)	34.14
	8B-Standard RL	88.37 (± 0.37)	38.06 (± 0.29)	81.10 (± 0.37)	11.03 (± 0.34)	42.78	35.18 (± 0.36)	37.95 (± 0.30)	31.30 (± 0.31)	34.81
8B-MENTOR	89.41* (± 0.56)	38.77* (± 0.60)	81.96* (± 0.52)	11.57* (± 0.42)	43.39	35.77* (± 0.36)	38.86* (± 0.56)	31.68* (± 0.36)	35.44	

Table 2: Main results comparing MENTOR against baselines across all evaluation benchmarks. The top table shows in-domain accuracy (%) on mathematical reasoning tasks. The bottom table shows out-of-domain performance on BFCL-v4 (accuracy %) and RAG (exact match %). Results are reported as mean (\pm standard deviation) over 10 runs. Overall scores are calculated differently: as a macro-average for MATH and RAG, and as an official weighted average for BFCL-v4. Values marked with an asterisk (*) denote statistical significance, determined by the Bonferroni pairwise test at a significance level of $\alpha < 0.05$.

tuning the student SLMs. For each model, we used the recommended sampling parameters from its official repository. The ‘search’ tool provided to the agent for the retrieval-based QA tasks is powered by a retrieval environment based on FlashRAG (Jin et al., 2025), using E5-base-v2 (Wang et al., 2022) as the retriever and the Dec. 2018 Wikipedia snapshot (Karpukhin et al., 2020) as the knowledge base. For retrieval-based tasks, we retrieve the top-5 results for each query. We provide further details on hyperparameters in Appendix A.5 and on prompts in Appendix B.

4.2 Experimental Results

The main results, presented in Table 2, demonstrate that our proposed framework, MENTOR, significantly outperforms the baselines.

Distillation Enables Effective Tool Use. On in-domain tasks, both SFT and our RL-based distillation clearly enable more effective tool use, leading to significant performance gains over the vanilla baselines. This confirms that transferring the teacher’s tool-calling ability is a broadly successful strategy for improving SLM performance.

RL-Based Distillation Achieves General Performance. On out-of-domain (OOD) benchmarks, the SFT baseline often shows minimal improvement or even performance degradation on OOD tasks, which we attribute to its tendency to overfit on the training domain. By contrast, RL-based approaches consistently yield significant performance gains in the OOD setting. This highlights the effectiveness of the RL-based distillation framework in instilling a more robust and generalizable problem-solving methodology, rather than merely encouraging the memorization of the teacher’s trajectories. For a more detailed analysis, the following section focuses on the performance of the Qwen2.5-7B model.

5 Analysis

5.1 Impact of RL-Based Distillation

RL Drives Effective Transfer. To evaluate how efficiently each model learns the teacher’s strategy, we analyze the alignment of its tool-use policy with the teacher’s patterns. As shown in Figure 3, there is a clear correlation between model performance and alignment score (AS). The SFT model clearly demonstrates imitation failure by learning a

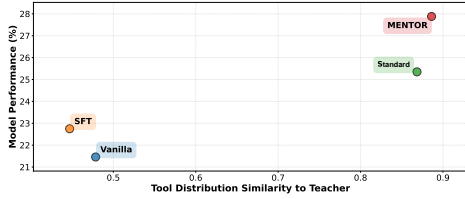


Figure 3: Correlation between task performance (Math) and alignment score (AS).

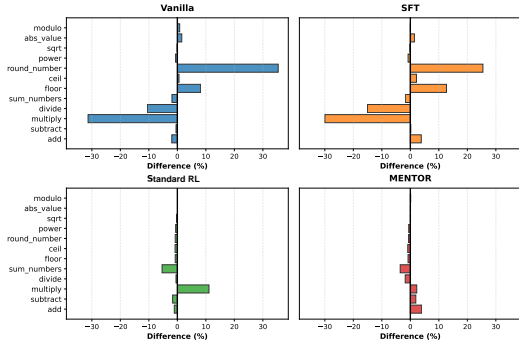


Figure 4: Comparison of tool invocation patterns between student models and the teacher. Each bar represents the percentage point difference in invocation frequency for a specific tool between the student model and the teacher.

policy that is the least similar to the teacher’s, confirming that it relies on superficial shortcuts rather than the intended methodology. In contrast, the RL-based methods achieve significantly higher performance and alignment with the teacher’s policy. Notably, MENTOR achieves the highest accuracy while learning a tool-use policy most similar to the expert teacher’s. This suggests that our RL-based distillation is the most effective method for transferring the efficient tool-use strategy.

RL Distillation Drives Policy Alignment. To evaluate the strategic alignment of each model with the teacher, Figure 4 compares their respective tool invocation patterns to the teacher’s reference policy, which is shown as the LLM’s tool usage in Figure 1. The baseline models, Vanilla and SFT, exhibit significant deviations from the teacher’s patterns, indicating that they learn a divergent and suboptimal tool-use policy. By contrast, the RL-based methods, such as standard RL and MENTOR, demonstrate a much closer alignment, with minimal differences across most tools. This provides strong visual evidence that our RL-based distillation is highly effective at teaching the SLM to internalize the teacher’s strategic methodology.

RL Enables Efficient Tool Use. Figure 5 shows the tool-use efficiency of each model by plotting the distribution of tool calls per question for both in-domain and out-of-domain tasks. While the

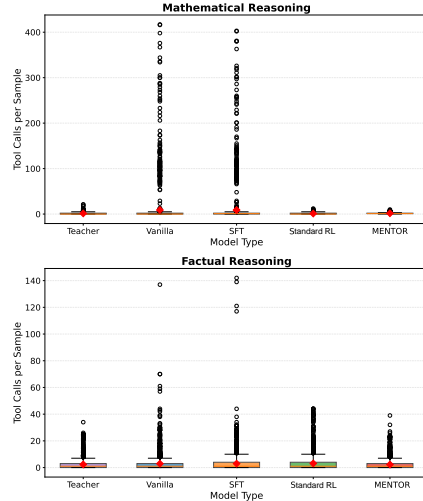


Figure 5: Tool-use efficiency on in-domain and out-of-domain tasks, measured by the distribution of tool calls per sample.

Teacher model is highly efficient, the Vanilla and SFT baselines tend to use tools inefficiently, evidenced by their wide distributions and numerous outliers. In contrast, our RL-based approach effectively transfers the teacher’s efficient strategy, maintaining a low and stable number of tool calls. It suggests that the RL-based approach successfully transfers the teacher’s tool-use efficiency and that this skill generalizes to out-of-domain tasks.

5.2 Impact of Teacher-Guided Reward

Reward Setting	Math	BFCL	RAG	AS
(1) R_c (Standard RL)	26.13	30.88	18.05	82.65
(2) $R_c + R_v$	27.80	30.96	8.35	80.94
(3) $R_c + R_{\text{Tool format}}$	26.61	28.95	7.83	79.73
(4) $R_c + R_{\text{Tool format}} + R_v$	26.76	30.60	9.42	83.54
(5) $R_c + R_a^{\text{strict}} + R_v$	26.88	30.85	18.05	85.25
(6) $R_c + R_a + R_v$ (Ours)	27.88	31.38	21.23	88.79

Table 3: Ablation study of the reward components. Performance is measured by the overall score on the Math, BFCL, and RAG benchmarks. AS represents the alignment score.

Ablation of Rewards We conduct an ablation study to isolate the contribution of each reward component by testing five distinct settings. The results are shown in Table 3. Specifically, **Setting (1)** serves the standard RL baseline, using only a reward for final answer correctness (R_a). **Setting (2)** introduces a tool validation reward (R_v) to penalize execution errors and invalid calls. **Setting (3)** adds a reward for the correct tool format, a strict structural constraint used in prior works (Qian et al., 2025). **Setting (4)** adds the validation reward to Setting (2). **Setting (5)** utilizes a strict alignment reward (R_a^{strict}) to enforce strict sequential alignment with the teacher’s trajectory. **Setting (6)** is MENTOR, which combines all rewards for the teacher-guided

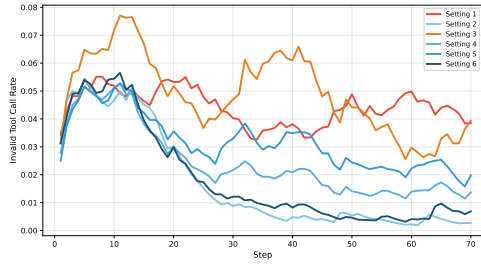


Figure 6: Invalid tool call rate over training steps

reward, employing a set-based alignment. The results show that our comprehensive reward design (Setting 6) is the best performer, achieving the highest scores across all task benchmarks and in policy alignment.

Referencing Teacher’s Trajectory Drives Performance. The results (Table 3) demonstrate a strong positive correlation (Pearson’s $r = 0.87$) between average performance and the alignment score (AS), indicating that the student’s success depends on how closely it references the teacher’s methodology. Notably, Settings 4, 5, and 6, which incorporate the teacher’s trajectory as a reference standard for reward calculation, achieve the highest alignment scores. This confirms that performance is enhanced by actively guiding the student to reference and adopt the teacher’s tool-use patterns.

Flexible Guidance Achieves Better Alignment than Strict Constraints. Our ablation study further reveals that the method of enforcing alignment is critical. We compare our flexible, set-based alignment reward (Setting 6) against a strict, sequence-based exact-match reward (Setting 5). The results (Table 3) demonstrate that strict sequence constraints are less effective than flexible guidance, as the latter achieves a higher Alignment Score and superior downstream performance. This indicates that strict constraints are excessively restrictive, punishing valid variations and hindering learning. In contrast, flexible guidance focuses on the strategic selection of tools, enabling the student to internalize the teacher’s methodology more effectively and achieve more robust alignment.

Validation Reward Reduces Invocation Errors. To demonstrate the effectiveness of our reward design in addressing tool invocation errors, Figure 6 presents the invalid tool call rate during training for each of the five reward settings from our ablation study. The results show that the choice of reward components leads to distinct learning behaviors. Settings that lack the tool validation reward (R_v),

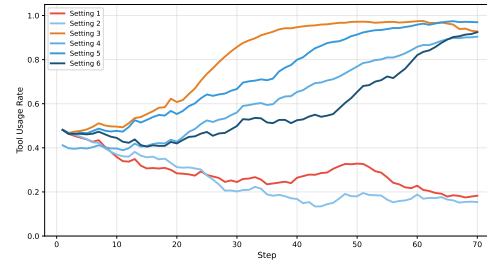


Figure 7: Tool usage rate over training steps

such as Settings 1 and 3, fail to effectively reduce their error rates, which remain high and unstable throughout training. In contrast, settings that include the tool validation reward (Settings 2, 4, 5, and 6) learn to avoid invalid calls, with their error rates dropping rapidly to near-zero. This suggests that directly penalizing invalid calls with R_v is a highly effective strategy for training a more reliable agent.

Teacher-Guided Rewards Foster Robust Tool Adoption. Figure 7 illustrates the impact of reward components on the model’s tendency to utilize tools. Without the guidance of the teacher’s trajectory (Settings 1 and 2), the student fails to recognize the value of tools and drifts toward a sub-optimal tool-avoidant strategy. Conversely, incorporating teacher alignment (Settings 3-6) actively encourages tool adoption. Crucially, while strict constraints (Setting 5) drive immediate but rigid usage, MENTOR (Setting 6) demonstrates a steadier learning curve that converges to high usage. This indicates that the student is not merely memorizing a simple heuristic to “use tools”, but is internalizing the teacher’s strategy, learning to deploy tools deliberately to solve complex problems.

6 Conclusion

To address the poor generalization of SFT and the inefficiency of simple outcome-reward RL in distilling tool-use into SLMs, we introduce MENTOR, a framework that combines reinforcement learning with a composite, teacher-guided reward. Extensive experiments demonstrate that MENTOR significantly improves cross-domain generalization by learning a policy that achieves both closer alignment with the teacher’s strategy and superior tool-use efficiency compared to baselines. Our ablation studies confirm that the synergistic combination of each component in our composite reward design is crucial for achieving this robust performance.

575 Limitations

576 While this study offers valuable insights, it is essen- 624
577 tial to acknowledge that several open challenges 625
578 remain. 626

579 **Extending to Other Models** Our current frame- 627
580 work focuses on the Qwen model series (Qwen2.5 628
581 and Qwen3), a choice guided by the design of our 629
582 RL-based approach. Our method is intended to re- 630
583 fine and generalize an existing tool-use policy and 631
584 thus performs best when the student model has 632
585 some initial ability. A key open challenge is adapt- 633
586 ing this framework to models that completely lack 634
587 this initial ability, as the exploratory feedback from 635
588 RL alone may be insufficient for them to learn ef- 636
589 fectively. A promising direction for future work 637
590 is to explore a two-stage, SFT-then-RL training 638
591 pipeline. Such an approach could first use an SFT 639
592 phase to instill a baseline policy before our RL- 640
593 based refinement is applied, thereby extending the 641
594 applicability of our method to a broader range of 642
595 models. 643

596 **Extending to Real-World Environments** Our 644
597 current work operates within a sandbox environ- 645
598 ment where tools are invoked via executable code. 646
599 A significant avenue for future research is to 647
600 extend our framework to operate in more com- 648
601 plex, tool-augmented environments, such as web 649
602 browsers, simulators, or desktop interfaces. In par- 650
603 ticular, integration with the Model Context Protocol 651
604 (MCP) (Anthropic, 2024)—which utilizes servers 652
605 that can respond to various scenarios—could sig- 653
606 nificantly enhance the capabilities of small agents 654
607 across a diverse range of real-world tasks. This 655
608 represents an important direction for future work. 656

609 Ethical Statements

610 This work contributes to the development of effi- 657
611 cient and capable artificial intelligence. By success- 658
612 fully distilling the complex tool-use capabilities of 659
613 large language models (LLMs) into smaller, more 660
614 efficient small language models (SLMs), MENTOR 661
615 accelerates the creation of functional, on-device 662
616 AI. This capability enables local deployment, re- 663
617 ducing reliance on expensive cloud infrastructure 664
618 and improving user privacy for agents that perform 665
619 complex tasks, such as mathematical reasoning and 666
620 information retrieval from external sources (includ- 667
621 ing the web). 668

622 However, the enhanced strategic competence and 669
623 tool-augmented abilities conferred by MENTOR 670

624 also introduce potential risks. Since our distilled 625
626 agents are capable of autonomous reasoning, web 627
628 retrieval, and code execution, they could be suscep- 628
629 tible to misuse. Potential malicious behaviors in- 629
630 clude the automated generation of harmful scripts, 630
631 the execution of unauthorized actions, or the spread 631
632 of misinformation via tool-based retrieval. To en- 632
633 sure responsible deployment, the integration of ro- 633
634 bust safeguards is essential. We emphasize that 634
635 addressing these ethical and safety concerns is an 635
636 important direction for future research and respon- 636

References 637

- 637 AI-MO. 2024. Aime. [https://huggingface.co/](https://huggingface.co/datasets/AI-MO/aimo-validation-aime) 638
639 [datasets/AI-MO/aimo-validation-aime](https://huggingface.co/datasets/AI-MO/aimo-validation-aime). 639
- 640 Mathematical Association of America AMC. 640
641 2023. 2023 AMC 12a and 12b: Amer- 641
642 ican mathematics competitions. [https:](https://www.maa.org/math-competitions/amc-12) 642
643 [//www.maa.org/math-competitions/amc-12](https://www.maa.org/math-competitions/amc-12). 643
644 Official competition information available at 644
645 the MAA website. Problem statements refer- 645
646 enced via the Art of Problem Solving archive: 646
647 [https://artofproblemsolving.com/wiki/](https://artofproblemsolving.com/wiki/index.php/2023_AMC_12A_Problems) 647
648 [index.php/2023_AMC_12A_Problems](https://artofproblemsolving.com/wiki/index.php/2023_AMC_12A_Problems) and 648
649 [https://artofproblemsolving.com/wiki/](https://artofproblemsolving.com/wiki/index.php/2023_AMC_12B_Problems) 649
650 [index.php/2023_AMC_12B_Problems](https://artofproblemsolving.com/wiki/index.php/2023_AMC_12B_Problems). Accessed: 650
651 2025-10-06. 651
- 652 Anthropic. 2024. Introducing the model context 652
653 protocol. [https://www.anthropic.com/news/](https://www.anthropic.com/news/model-context-protocol/) 653
654 [model-context-protocol/](https://www.anthropic.com/news/model-context-protocol/). 654
- 655 Mingyang Chen, Linzhuang Sun, Tianpeng Li, Haoze 655
656 Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang, 656
657 Jeff Z Pan, Wen Zhang, Huajun Chen, and 1 oth- 657
658 ers. 2025a. Learning to reason with search for 658
659 llms via reinforcement learning. *arXiv preprint* 659
660 *arXiv:2503.19470*. 660
- 661 Yang Chen, Zhuolin Yang, Zihan Liu, Chankyu Lee, 661
662 Peng Xu, Mohammad Shoeybi, Bryan Catanzaro, 662
663 and Wei Ping. 2025b. Acereason-nemotron: Advanc- 663
664 ing math and code reasoning through reinforcement 664
665 learning. *arXiv preprint arXiv:2505.16400*. 665
- 666 Li Chenglin, Qianglong Chen, Liangyue Li, Caiyu 666
667 Wang, Feng Tao, Yicheng Li, Zulong Chen, and Yin 667
668 Zhang. 2024. [Mixed distillation helps smaller lan- 668](#)
669 [guage models reason better](#). In *Findings of the Asso- 669*
670 *ciation for Computational Linguistics: EMNLP 2024*, 670
671 pages 1673–1690, Miami, Florida, USA. Association 671
672 for Computational Linguistics. 672
- 673 Chengwei Dai, Kun Li, Wei Zhou, and Songlin Hu. 673
674 2024a. Beyond imitation: Learning key reasoning 674
675 steps from dual chain-of-thoughts in reasoning distil- 675
676 lation. *arXiv preprint arXiv:2405.19737*. 676

677	Chengwei Dai, Kun Li, Wei Zhou, and Songlin Hu.	Chengsong Huang, Wenhao Yu, Xiaoyang Wang, Hong-	733
678	2024b. Improve student’s reasoning generalizabil-	ming Zhang, Zongxia Li, Ruosen Li, Jiaxin Huang,	734
679	ity through cascading decomposed cots distillation.	Haitao Mi, and Dong Yu. 2025. R-zero: Self-	735
680	<i>arXiv preprint arXiv:2405.19842</i> .	evolving reasoning llm from zero data. <i>arXiv</i>	736
		<i>preprint arXiv:2508.05004</i> .	737
681	Chengwei Dai, Kun Li, Wei Zhou, and Songlin Hu.	Tenghao Huang, Dongwon Jung, Vaibhav Kumar, Mo-	738
682	2025. Capture the key in reasoning to enhance CoT	hammad Kachuee, Xiang Li, Puyang Xu, and Muhao	739
683	distillation generalization . In <i>Proceedings of the 63rd</i>	Chen. 2024. Planning and editing what you retrieve	740
684	<i>Annual Meeting of the Association for Computational</i>	for enhanced tool learning . In <i>Findings of the Associ-</i>	741
685	<i>Linguistics (Volume 1: Long Papers)</i> , pages 441–465.	<i>ation for Computational Linguistics: NAACL 2024</i> ,	742
686	Association for Computational Linguistics.	pages 975–988, Mexico City, Mexico. Association	743
		for Computational Linguistics.	744
687	Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang,	Tatsuro Inaba, Hirokazu Kiyomaru, Fei Cheng, and	745
688	Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin	Sadao Kurohashi. 2023. MultiTool-CoT: GPT-3	746
689	Chi, and Wanjun Zhong. 2025. Retool: Reinforce-	can use multiple external tools with chain of thought	747
690	ment learning for strategic tool use in llms. <i>arXiv</i>	prompting . In <i>Proceedings of the 61st Annual Meet-</i>	748
691	<i>preprint arXiv:2504.11536</i> .	<i>ing of the Association for Computational Linguistics</i>	749
		<i>(Volume 2: Short Papers)</i> , pages 1522–1532, Toronto,	750
692	Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo	Canada. Association for Computational Linguistics.	751
693	Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang		
694	Chen, Runxin Xu, and 1 others. 2024. Omni-	Jiajie Jin, Yutao Zhu, Zhicheng Dou, Guanting Dong,	752
695	math: A universal olympiad level mathematic bench-	Xinyu Yang, Chenghao Zhang, Tong Zhao, Zhao	753
696	mark for large language models. <i>arXiv preprint</i>	Yang, and Ji-Rong Wen. 2025. Flashrag: A modular	754
697	<i>arXiv:2410.07985</i> .	toolkit for efficient retrieval-augmented generation	755
		research . In <i>Companion Proceedings of the ACM</i>	756
698	Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon,	<i>on Web Conference 2025, WWW 2025, Sydney, NSW,</i>	757
699	Pengfei Liu, Yiming Yang, Jamie Callan, and Graham	<i>Australia, 28 April 2025 - 2 May 2025</i> , pages 737–	758
700	Neubig. 2023. Pal: Program-aided language models.	740. ACM.	759
701	In <i>International Conference on Machine Learning</i> ,		
702	pages 10764–10799. PMLR.	Leslie Pack Kaelbling, Michael L Littman, and An-	760
		drew W Moore. 1996. Reinforcement learning: A	761
703	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen,	survey. <i>Journal of artificial intelligence research</i> ,	762
704	Yujia Yang, Minlie Huang, Nan Duan, and Weizhu	4:237–285.	763
705	Chen. 2024. Tora: A tool-integrated reasoning agent		
706	for mathematical problem solving. <i>The Twelfth Inter-</i>	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric	764
707	<i>national Conference on Learning Representations</i> .	Wallace, and Colin Raffel. 2023. Large language	765
		models struggle to learn long-tail knowledge. In	766
708	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	<i>International conference on machine learning</i> , pages	767
709	Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-	15696–15707. PMLR.	768
710	cob Steinhardt. 2021. Measuring mathematical prob-		
711	lem solving with the MATH dataset . In <i>Proceedings</i>	Minki Kang, Jongwon Jeong, Seanie Lee, Jaewoong	769
712	<i>of the Neural Information Processing Systems Track</i>	Cho, and Sung Ju Hwang. 2025. Distilling llm agent	770
713	<i>on Datasets and Benchmarks 1, NeurIPS Datasets</i>	into small models with retrieval and code tools. <i>arXiv</i>	771
714	<i>and Benchmarks 2021, December 2021, virtual</i> .	<i>preprint arXiv:2505.17612</i> .	772
715	Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023.	Vladimir Karpukhin, Barlas Oguz, Sewon Min,	773
716	Large language models are reasoning teachers . In	Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi	774
717	<i>Proceedings of the 61st Annual Meeting of the As-</i>	Chen, and Wen-tau Yih. 2020. Dense passage re-	775
718	<i>sociation for Computational Linguistics (Volume 1:</i>	trieval for open-domain question answering. In	776
719	<i>Long Papers)</i> , pages 14852–14882.	<i>EMNLP (1)</i> , pages 6769–6781.	777
720	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara,	Yilun Kong, Jingqing Ruan, Yihong Chen, Bin Zhang,	778
721	and Akiko Aizawa. 2020. Constructing A multi-hop	Tianpeng Bao, Shiwei Shi, Guoqing Du, Xiaoru Hu,	779
722	QA dataset for comprehensive evaluation of reason-	Hangyu Mao, Ziyue Li, and 1 others. 2023. Tptu-	780
723	ing steps . In <i>Proceedings of the 28th International</i>	v2: Boosting task planning and tool usage of large	781
724	<i>Conference on Computational Linguistics, COLING</i>	language model-based agents in real-world systems.	782
725	<i>2020, Barcelona, Spain (Online), December 8-13,</i>	<i>arXiv preprint arXiv:2311.11315</i> .	783
726	2020, pages 6609–6625. International Committee on		
727	Computational Linguistics.	Aitor Lewkowycz, Anders Andreassen, David Dohan,	784
		Ethan Dyer, Henryk Michalewski, Vinay Ramasesh,	785
728	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	Ambrose Slone, Cem Anil, Imanol Schlag, Theo	786
729	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,	Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy	787
730	and Weizhu Chen. 2021. Lora: Low-rank adap-	Gur-Ari, and Vedant Misra. 2022. Solving quan-	788
731	tation of large language models. <i>arXiv preprint</i>	titative reasoning problems with language models .	789
732	<i>arXiv:2106.09685</i> .	<i>Preprint</i> , arXiv:2206.14858.	790

791	Dacheng Li, Shiyi Cao, Tyler Griggs, Shu Liu, Xi- angxi Mo, Eric Tang, Sumanth Hegde, Kourosh 792 Hakhmaneshi, Shishir G Patil, Matei Zaharia, and 793 1 others. 2025. Llms can easily learn to reason from 794 demonstrations structure, not content, is what mat- 795 ters! <i>arXiv preprint arXiv:2502.07374</i> .	847
796		848
797	Huanxuan Liao, Shizhu He, Yao Xu, Yuanzhe Zhang, 798 Kang Liu, and Jun Zhao. 2025. Neural-symbolic 799 collaborative distillation: Advancing small language 800 models for complex reasoning tasks. In <i>Proceedings 801 of the AAAI Conference on Artificial Intelligence</i> , 802 volume 39, pages 24567–24575.	849
803		850
804	Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen 805 Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong, 806 Ernie Chang, Yangyang Shi, Raghuraman Krish- 807 namoorthi, and 1 others. 2024. Mobilellm: Opti- 808 mizing sub-billion parameter language models for 809 on-device use cases. In <i>Forty-first International Con- ference on Machine Learning</i> .	851
810		852
811	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jian- 812 guang Lou, Chongyang Tao, Xiubo Geng, Qingwei 813 Lin, Shifeng Chen, and Dongmei Zhang. 2025a. Wiz- 814 ardmath: Empowering mathematical reasoning for 815 large language models via reinforced evol-instruct. 816 <i>The Twelfth International Conference on Learning Representations</i> .	853
817		854
818	Ne Luo, Aryo Pradipta Gema, Xuanli He, Emile 819 Van Krieken, Pietro Lesci, and Pasquale Minervini. 820 2025b. Self-training large language models for 821 tool-use without demonstrations. <i>arXiv preprint arXiv:2502.05867</i> .	855
822		856
823	Yuanjie Lyu, Chengyu Wang, Jun Huang, and Tong 824 Xu. 2025. From correction to mastery: Reinforced 825 distillation of large language model agents. <i>arXiv preprint arXiv:2509.14257</i> .	857
826		858
827	Bhargavi Paranjape, Scott Lundberg, Sameer Singh, 828 Hannaneh Hajishirzi, Luke Zettlemoyer, and 829 Marco Tulio Ribeiro. 2023. Art: Automatic 830 multi-step reasoning and tool-use for large language models. <i>arXiv preprint arXiv:2303.09014</i> .	859
831		860
832	Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, 833 Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph 834 E. Gonzalez. 2025. The berkeley function calling 835 leaderboard (bfc): From tool use to agentic eval- 836 uation of large language models. In <i>Forty-second International Conference on Machine Learning</i> .	861
837		862
838	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, 839 Noah Smith, and Mike Lewis. 2023. Measuring and 840 narrowing the compositionality gap in language mod- 841 els . In <i>Findings of the Association for Computational 842 Linguistics: EMNLP 2023</i> , pages 5687–5711, Singa- 843 pore. Association for Computational Linguistics.	863
844		864
845	Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, 846 Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. 2025. Toolrl: Reward is all tool learning needs. <i>arXiv preprint arXiv:2504.13958</i> .	865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901

902		Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	959
903		Shafran, Karthik Narasimhan, and Yuan Cao. 2023.	960
904		React: Synergizing reasoning and acting in language	961
		models. In <i>International Conference on Learning</i>	962
		<i>Representations (ICLR)</i> .	963
905	Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun	Maxwell J Yin, Dingyi Jiang, Yongbing Chen, Boyu	964
906	Luo, Weikang Shi, Renrui Zhang, Linqi Song,	Wang, and Charles Ling. 2025. Enhancing general-	965
907	Mingjie Zhan, and Hongsheng Li. 2024a. Mathcoder:	ization in chain of thought reasoning for smaller	966
908	Seamless code integration in llms for enhanced math-	models. <i>arXiv preprint arXiv:2501.09804</i> .	967
909	ematical reasoning. In <i>The Twelfth International</i>		
910	<i>Conference on Learning Representations</i> .		
911	Liang Wang, Nan Yang, Xiaolong Huang, Binx-	Yuanqing Yu, Zhefan Wang, Weizhi Ma, Shuai Wang,	968
912	ing Jiao, Linjun Yang, Daxin Jiang, Rangan Ma-	Chuhan Wu, Zhiqiang Guo, and Min Zhang. 2024.	969
913	jumder, and Furu Wei. 2022. Text embeddings by	Steptool: Enhancing multi-step tool usage in llms	970
914	weakly-supervised contrastive pre-training. <i>CoRR</i> ,	through step-grained reinforcement learning. <i>arXiv</i>	971
915	abs/2212.03533.	<i>preprint arXiv:2410.07745</i> .	972
916	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang,	Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun	973
917	Yunzhu Li, Hao Peng, and Heng Ji. 2024b. Exe-	Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song,	974
918	cutable code actions elicit better llm agents. In <i>Forty-</i>	Mingjie Zhan, and 1 others. 2023. Solving chal-	975
919	<i>first International Conference on Machine Learning</i> .	lenging math word problems using gpt-4 code in-	976
920	Zengzhi Wang, Fan Zhou, Xuefeng Li, and Pengfei	terpreter with code-based self-verification. <i>arXiv</i>	977
921	Liu. 2025. Octothinker: Mid-training incentivizes	<i>preprint arXiv:2308.07921</i> .	978
922	reinforcement learning scaling. <i>arXiv preprint</i>		
923	<i>arXiv:2506.20512</i> .		
924	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,		
925	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,		
926	Maarten Bosma, Denny Zhou, Donald Metzler, and		
927	1 others. 2022. Emergent abilities of large language		
928	models. <i>arXiv preprint arXiv:2206.07682</i> .		
929	Haoze Wu, Yunzhi Yao, Wenhao Yu, Huajun Chen,		
930	and Ningyu Zhang. 2025a. Recode: Updating code		
931	api knowledge with reinforcement learning. <i>arXiv</i>		
932	<i>preprint arXiv:2506.20495</i> .		
933	Yongliang Wu, Yizhou Zhou, Zhou Ziheng, Yingzhe		
934	Peng, Xinyu Ye, Xinting Hu, Wenbo Zhu, Lu Qi,		
935	Ming-Hsuan Yang, and Xu Yang. 2025b. On the		
936	generalization of sft: A reinforcement learning per-		
937	spective with reward rectification. <i>arXiv preprint</i>		
938	<i>arXiv:2508.05629</i> .		
939	Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu,		
940	Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Su-		
941	jian Li. 2024. Watch every step! Llm agent learning		
942	via iterative step-level process refinement. <i>arXiv</i>		
943	<i>preprint arXiv:2406.11176</i> .		
944	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,		
945	Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-		
946	hong Tu, Jingren Zhou, Junyang Lin, Keming Lu,		
947	Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang		
948	Ren, and Zhenru Zhang. 2024. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement . <i>ArXiv</i> .		
949			
950			
951	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,		
952	William Cohen, Ruslan Salakhutdinov, and Christo-		
953	pher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering .		
954	In <i>Proceedings of the 2018 Conference on Empirical</i>		
955	<i>Methods in Natural Language Processing</i> , pages		
956	2369–2380, Brussels, Belgium. Association for Com-		
957	putational Linguistics.		
958			

A Details for Experimental Setup

A.1 Training Dataset Details

To train our models for effective tool use, we construct a high-quality dataset of reference trajectories, as described in Section 3.1. To this end, we use the nvidia/AceReason-Math² (Chen et al., 2025b) dataset as our source, leveraging its verified questions and ground-truth answers to ensure reliability.

The construction process involves providing these questions to our teacher model, Qwen3-235B, which generated solution trajectories using a calculator tool. We then filter these outputs, retaining only the **successful trajectories** where the teacher’s final answer matched the ground-truth from the source dataset. This verification process yields a final training set of 1.27k single- and multi-turn trajectories that exemplify successful tool use.

A.2 Versions of Models

We employ Qwen3-235B-Thinking (Qwen3, 2025) as the teacher model and four student models from both Qwen3 (Qwen3, 2025) and Qwen2.5 (Yang et al., 2024) families. The specific model versions and their Hugging Face identifiers are listed in Table 4.

Role	Model	Hugging Face Identifier
Teacher	Qwen3-235B-Thinking	Qwen/Qwen3-235B-A22B-Thinking-2507-FP8
Student	Qwen2.5-1.5B-Instruct	Qwen/Qwen2.5-1.5B-Instruct
	Qwen2.5-7B-Instruct	Qwen/Qwen2.5-7B-Instruct
	Qwen3-1.7B	Qwen/Qwen3-1.7B
	Qwen3-8B	Qwen/Qwen3-8B

Table 4: Model versions used in experiments

A.3 Benchmarks

We use the MATH³, Omni-MATH-512⁴, AIME24⁵, AIME25⁶, amc23⁷, and minervamath⁸ datasets. All datasets are publicly available for research use.

²<https://hf.co/datasets/nvidia/AceReason-Math>

³<https://hf.co/datasets/prithivMLmods/Math-Forge-Hard>

⁴<https://hf.co/datasets/Heng1999/Omni-MATH-512>

⁵<https://hf.co/datasets/math-ai/aime24>

⁶<https://hf.co/datasets/math-ai/aime25>

⁷<https://hf.co/datasets/math-ai/amc23>

⁸<https://hf.co/datasets/math-ai/minervamath>

A.4 Teacher-Student Alignment Metric

To quantify alignment in tool usage patterns, we compare the distributions of tool calls across the 12 available tools, which are defined in Table 11. Given the teacher’s tool distribution P and a student’s distribution Q , we use an alignment score based on the Jensen-Shannon Divergence (JSD), defined as:

$$\text{Alignment}(P, Q) = 1 - \text{JSD}(P, Q)$$

$$\text{where } \text{JSD}(P, Q) = \sqrt{\frac{1}{2}\mathbb{D}_{KL}(P\|M) + \frac{1}{2}\mathbb{D}_{KL}(Q\|M)}$$

$$\text{and } M = \frac{1}{2}(P + Q)$$
(8)

The score is bounded between 0 and 1, where 1 signifies a perfect match. \mathbb{D}_{KL} denotes the Kullback-Leibler divergence.

A.5 Hyperparameters for Training and Inference

Our training is conducted on 1×8 Nvidia H100 80G GPUs, with full parameter optimization and gradient checkpointing. We provide some important parameter settings in Tables 5 and 6.

Parameter	Value
Learning Rate	1e-6
Optimizer	AdamW
Epochs	2
Train Batch Size	16
Mini-batch Size	8
Max Sequence Length	32768
Max Response Length	8192
Number of Rollout	10
Tensor Model Parallel Size	2
Rollout Temperature (Qwen3)	0.6
Rollout Temperature (Qwen2.5)	0.7
GPU Utilization Ratio	0.8
KL Loss Coefficient	0.001
Clip Ratio	0.2
Reward Weights (w_c, w_a, w_v)	0.7, 0.3, 0.1

Table 5: Implementation details of MENTOR.

Parameter	Value
Learning Rate	1e-7
Optimizer	AdamW
Epochs	1
Train Batch Size	4
Gradient Accumulation Steps	16
Weight decay	0.033
Gradient Accum	8

Table 6: Implementation details of SFT.

B Details of Instruction Prompts

We utilize prompts for both reference trajectory generation and model evaluation. The specific

1028 prompt for reference trajectory generation is shown
1029 in Table 8. For evaluation, we assess the robustness
1030 of MENTOR and the baseline models using prompts
1031 tailored to each of the three domains. All evalua-
1032 tion prompts for the Qwen2.5 and Qwen3 models
1033 are created by adapting the official chat templates
1034 provided with each model’s tokenizer. The prompts
1035 are shown in Tables 7, 8, 9, and 10. The in-
1036 struction prompts used for evaluating the BFCL-v4
1037 benchmark directly follow the format and content
1038 specified on the official benchmark website⁹, en-
1039 suring consistency with prior work¹⁰.

1040 C Tool Execution via Remote Server

1041 All tools described in Tables 11, and 12 are imple-
1042 mented as executable Python code. For tool exe-
1043 cution, we use a FastAPI-based remote execution
1044 server, following the base architecture of the ReCall
1045 framework (Chen et al., 2025a). Our implementa-
1046 tion code is publicly available¹¹.

1047 When an agent invokes a tool, the system sends
1048 the tool’s Python code to the remote server via
1049 HTTP API. The server executes the code in a pre-
1050 configured environment with necessary libraries
1051 installed and returns the result.

⁹https://gorilla.cs.berkeley.edu/blogs/15_bfcl_v4_web_search.html

¹⁰<https://github.com/ShishirPatil/gorilla/tree/main/berkeley-function-call-leaderboard>

¹¹<https://anonymous.4open.science/r/MENTOR-F6E7/>

Prompt for Qwen2.5 + MATH

SYSTEM:

system

You are Qwen, created by Alibaba Cloud. You are a helpful assistant.

Tools

You may call one or more functions to assist with the user query.

You are provided with function signatures within `<tools></tools>` XML tags:

`<tools>`

```
{"type": "function", "function": {"name": "add", ...}}
```

```
{"type": "function", "function": {"name": "subtract", ...}}
```

```
{"type": "function", "function": {"name": "multiply", ...}}
```

...

```
{"type": "function", "function": {"name": "modulo", ...}}
```

`</tools>`

For each function call, return a json object with function name and arguments within `<tool_call></tool_call>` XML tags:

`<tool_call>`

```
{"name": <function-name>, "arguments": <args-json-object>}
```

`</tool_call>`

USER:

Question: Cities A and B are 45 miles apart. Alicia lives in A and Beth lives in B . Alicia bikes towards B at 18 miles per hour. Leaving at the same time, Beth bikes toward A at 12 miles per hour. How many miles from City A will they be when they meet?

If you have got the answer, enclose it within `\boxed{}` with latex format.

Table 7: Prompt for Qwen2.5 + MATH

Prompt for Qwen3 + MATH

SYSTEM:

system
Tools

You may call one or more functions to assist with the user query.

You are provided with function signatures within `<tools></tools>` XML tags:

```
<tools>
{"type": "function", "function": {"name": "add", ...}}
{"type": "function", "function": {"name": "subtract", ...}}
{"type": "function", "function": {"name": "multiply", ...}}
...
{"type": "function", "function": {"name": "modulo", ...}}
</tools>
```

For each function call, return a json object with function name and arguments within `<tool_call></tool_call>` XML tags:

```
<tool_call>
{"name": <function-name>, "arguments": <args-json-object>}
</tool_call>
```

USER:

Question: Cities A and B are 45 miles apart. Alicia lives in A and Beth lives in B . Alicia bikes towards B at 18 miles per hour. Leaving at the same time, Beth bikes toward A at 12 miles per hour. How many miles from City A will they be when they meet?

If you have got the answer, enclose it within `\boxed{}` with latex format.

Table 8: Prompt for Qwen3 + MATH

Prompt for Qwen2.5 + RAG

SYSTEM:

system

You are Qwen, created by Alibaba Cloud. You are a helpful assistant.

Tools

You may call one or more functions to assist with the user query.

You are provided with function signatures within `<tools></tools>` XML tags:

```
<tools>
{"type": "function", "function": {
  "name": "wikipedia_search",
  "description": "Search Wikipedia for a given query.",
  "parameters": {
    "type": "object",
    "properties": {
      "query": {
        "type": "string",
        "description": "Query to search for."
      },
      "top_n": {
        "type": "integer",
        "description": "Number of results to return. The default value is 5.",
        "default": 5
      }
    },
    "required": ["query"]
  }
}
</tools>
```

For each function call, return a json object with function name and arguments within `<tool_call></tool_call>` XML tags:

```
<tool_call>
{"name": <function-name>, "arguments": <args-json-object>}
</tool_call>
```

USER:

Question: What is the capital of France?

If you have got the answer, enclose it within `\boxed{}` with latex format.

Table 9: Prompt for Qwen2.5 + RAG

Prompt for Qwen3 + RAG

SYSTEM:

system
Tools

You may call one or more functions to assist with the user query.

You are provided with function signatures within `<tools></tools>` XML tags:

```
<tools>
{"type": "function", "function": {
  "name": "wikipedia_search",
  "description": "Search Wikipedia for a given query.",
  "parameters": {"type": "object", "properties": {
    "query": {"type": "string", "description": "Query to search for."},
    "top_n": {"type": "integer", "description": "Number of results to
      return. The default value is 5.", "default": 5}},
    "required": ["query"]}]}
</tools>
```

For each function call, return a json object with function name and arguments within `<tool_call></tool_call>` XML tags:

```
<tool_call>
{"name": <function-name>, "arguments": <args-json-object>}
</tool_call>
```

USER:

Question: What is the capital of France?

If you have got the answer, enclose it within `\boxed{}` with latex format.

Table 10: Prompt for Qwen3 + RAG

Function	Description	Parameter Name	Parameter Description
add	Add two numbers together	firstNumber secondNumber	The first number The second number
subtract	Subtract one number from another	minuend subtrahend	The number to subtract from The number to subtract
multiply	Multiply two numbers together	firstNumber secondNumber	The first number The second number
divide	Divide one number by another	numerator denominator	The number to be divided The number to divide by
sum_numbers	Calculate the sum of an array of numbers	numbers	Array of numbers to sum
floor	Calculate the floor of a number	number	Number to find the floor of
ceil	Calculate the ceil of a number	number	Number to find the ceil of
round_number	Round a number to the nearest integer	number	Number to round
power	Calculate base raised to the power of exponent	base exponent	The base number The exponent
sqrt	Calculate the square root of a number	number	Number to find the square root of
abs_value	Calculate the absolute value of a number	number	Number to find the absolute value of
modulo	Calculate the modulo of two numbers	dividend divisor	The dividend The divisor

Table 11: Math Tool Functions

Function	Description	Parameter Name (Type)	Parameter Description
wikipedia_search	Search Wikipedia for a given query.	query (string) top_n (integer)	Query to search for. Number of results to return. (Optional, default: 5)

Table 12: RAG Tool Function