



PORTool: Tool-Use LLM Training with Rewarded Tree

Feijie Wu¹, Weiwu Zhu², Yuxiang Zhang², Soumya Chatterjee², Jiarong Zhu², Fan Mo², Rodin Luo², Jing Gao¹

¹Purdue University, ²Apple

Work done during an internship at Apple.

Current tool-use large language models (LLMs) are trained on static datasets, enabling them to interact with external tools and perform multi-step, tool-integrated reasoning, which produces tool-call trajectories. However, these models imitate how a query is resolved in a generic tool-call routine, thereby failing to explore possible solutions and demonstrating limited performance in an evolved, dynamic tool-call environment. In this work, we propose PORTool, a reinforcement learning (RL) method that encourages a tool-use LLM to explore various trajectories yielding the correct answer. Specifically, this method starts with generating multiple rollouts for a given query, and some of them share the first few tool-call steps, thereby forming a tree-like structure. Next, we assign rewards to each step, based on its ability to produce a correct answer and make successful tool calls. A shared step across different trajectories receives the same reward, while different steps under the same fork receive different rewards. Finally, these step-wise rewards are used to calculate fork-relative advantages, blended with trajectory-relative advantages, to train the LLM for tool use. The experiments utilize 17 tools to address user queries, covering both time-sensitive and time-invariant topics. We conduct ablation studies to systematically justify the necessity and the design robustness of step-wise rewards. Furthermore, we compare the proposed PORTool with other training approaches and demonstrate significant improvements in final accuracy and the number of tool-call steps.

Date: October 31, 2025

1 Introduction

As a representative paradigm for synergizing reasoning and action, the ReAct framework (Yao et al., 2023b) activates large language models (LLMs) to combine reasoning (step-by-step analysis) and acting (tool calls), achieving strong problem-solving performance. Early works (Yao et al., 2023a,b; Wei et al., 2022; Schick et al., 2023; Wu et al., 2025) enable this capability primarily through prompt engineering, in which carefully crafted demonstrations guide the model to alternate between generating reasoning traces and invoking tools. However, these approaches provide reasoning within a generic context rather than learning from actual interactions with external tools, thereby limiting their effectiveness in solving complex, multi-step tasks. This limitation highlights the importance of exploring *how to augment LLMs’ reasoning and acting capabilities through interactive tool use*.

Reinforcement learning (RL) has emerged as a promising direction to address this gap, leading to a growing number of research works (Qian et al., 2025; Zhang et al., 2025; Dong et al., 2025; Zeng et al., 2025; Jin et al., 2025; Li et al., 2025; Singh et al., 2025; Xue et al., 2025; Feng et al., 2025a; Wei et al., 2025; Lin and Xu, 2025). Most existing works concentrate on a narrow set of tools, such as web search or program compilers, which restricts the trained model’s applicability to a wide range of function-like tools with rigid formatting requirements and specialized functionalities. ToolRL (Qian et al., 2025) and Tool-N1 (Zhang et al., 2025) broaden the scope by training LLMs to interact with diverse tool schemas. Their datasets consist of structured inputs, including queries, tool descriptions, and, if available, historical tool-call trajectories, paired with the expected tool calls for subsequent steps. Both methods adopt group-relative policy optimization (GRPO) (Shao et al., 2024) for fine-tuning, where the model generates various tool-call steps for each input and receives rewards based on consistency with the expected tool calls, thereby reinforcing tool-use capability.

These two approaches, however, face several limitations because they assume the existence of ground-truth tool calls at every intermediate step. First, constructing such datasets is economically impractical in real-world settings, as it requires substantial human effort to understand tool execution mechanisms and annotate each step. Second, this design encourages the model to imitate predefined tool calls. While effective for producing generic routines that solve common queries, it diminishes the model’s ability to explore alternative trajectories. As a result, potentially valid solutions are either ignored or penalized for diverging from the expected calls. This limits the model’s ability to handle real-time queries, as illustrated in Figure 1. In this example, a model trained by ToolRL (Qian et al., 2025) attempts to finish the task by following a generic tool-call routine. However, the tool raises an error, leading the model to generate a no-answer response without considering whether it misunderstood the user’s intent.

In this work, we begin by collecting a set of real user queries that require both factual and real-time information, along with a suite of executable tools. Instead of providing a generic tool-call routine, we build an agentic system that consists of a tool-use agent and an evaluation agent: The tool-use agent enables the LLM to interact with external tools to generate complete tool-call trajectories, while the evaluation agent assesses the correctness of each trajectory based on its final answer to the input query. With this system, we follow the existing tool-use training works (Jin et al., 2025; Qu et al., 2025; Dong et al., 2025) to reinforce the LLM on the tool-use agent through GRPO. Specifically, the model generates multiple tool-call trajectories independently, and GRPO computes the advantages across the trajectory level, assigning identical advantages to all tokens within each trajectory.

Apparently, GRPO neglects the varying contributions of each step to the final result. In particular, even an incorrect trajectory may contain informative intermediate steps that, if recombined differently, could bring a correct answer. However, such potential is lost when all steps are assigned with uniform rewards (Zeng et al., 2025). Conversely, certain steps within a correct trajectory may have no influence on the final result. These observations emphasize the need to *quantify each step’s importance with the help of the evaluation agent*.

Toward the goal, we propose PORTool, a **P**olicy **O**ptimization algorithm with **R**ewarded Tree for **tool**-use LLMs. PORTool introduces *tree rollout*, in which certain tool-call steps branch into multiple trajectories and eventually lead to different subsequent tool calls. After evaluating each trajectory, the evaluation agent assigns a reward to each step that considers both the correctness of the final answer and its formatting compliance. Correctness is weighted more strongly in the reward signals, encouraging the model to favor a step that yields the correct answer, even though its format is incorrect. Notably, a step shared by various trajectories is assigned an identical reward, while different steps under the same forks receive distinct rewards. Finally, with these step-wise rewards, the tool-use agent computes fork-relative advantages that reflect the most effective steps within each fork. These are combined with trajectory-relative advantages that characterize the holistic performance of each trajectory to train the LLM, enabling it to generate more effective tool-call steps. Figure 1 illustrates the superiority of the proposed PORTool, where the trained model accurately understands the user’s intent, takes appropriate step-by-step tool calls, and provides a reliable response.

Contributions. Throughout this work, our contributions are highlighted as follows:

- To the best of our knowledge, this is the first work to reinforce the LLM for tool use by interacting with numerous executable tools, demonstrating evolved and dynamic outputs based on real-time information.

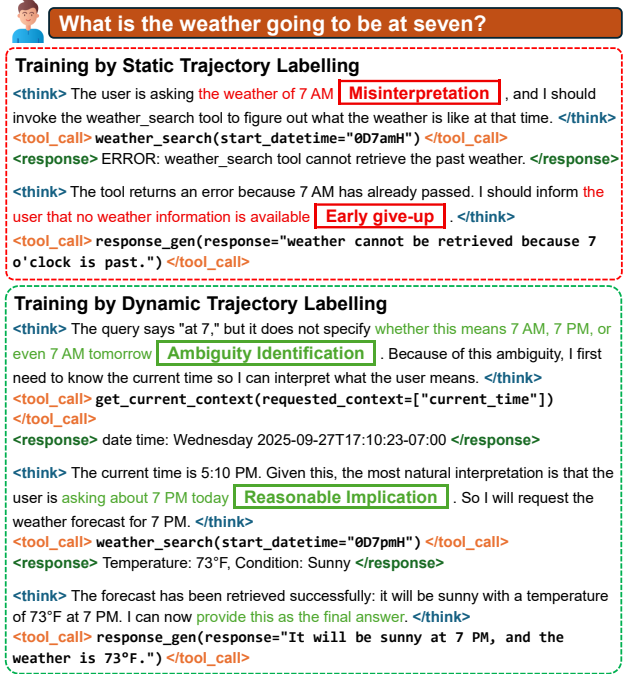


Figure 1 Training by labeled static tool-call trajectories cannot handle a real-time query. The example is generated from ToolRL (Qian et al., 2025).

- We propose an algorithm, PORTool, that integrates an innovative tree rollout strategy and designs a reward function that quantifies the importance of each step. Moreover, we theoretically figure out an optimal coefficient setting for combining step-relative and trajectory-relative advantages while reinforcing an LLM for tool use.
- We conduct comprehensive ablation studies to validate both the necessity and effectiveness of the proposed reward function design. In addition, we compare the proposed PORTool with other policy optimization algorithms and demonstrate superiority in terms of final accuracy and the number of tool-call steps.

2 Related Works

RL algorithms for enabling LLMs with reasoning. The rapid development of LLMs highlights the pivotal role of RL algorithms in training (Schulman et al., 2017; Rafailov et al., 2023; Ahmadian et al., 2024; Yu et al., 2025; Wu et al., 2024). Recent breakthroughs such as DeepSeek-R1 (Guo et al., 2025) disclose that RL algorithms can endow LLMs with chain-of-thought reasoning abilities (Wei et al., 2022) by assigning verifiable or rule-based rewards (Zhu et al., 2025; Yeo et al., 2025; Havrilla et al., 2024). This insight has spurred a wave of new policy optimization approaches (Yu et al., 2025; Zheng et al., 2025; Ahmadian et al., 2024; Liu et al., 2025; Chu et al., 2025b). In multi-step tool-use scenarios, these approaches assign a single reward to the entire trajectory, which obscures the varying significance of individual steps and fails to recognize how intermediate tool calls shape the final outcome. To tackle the limitation, GiGPO (Feng et al., 2025b) assigns rewards at the step level and normalizes them across actions within the same state when training the policy model. However, this framework does not extend to our scenario, as the state information is not explicitly defined.

Among those RL works, the most related concurrent work is ARPO (Dong et al., 2025). ARPO proposes an entropy-based adaptive rollout mechanism that generates multiple rollouts forming a tree structure. Specifically, forks are created when a significant probability deviation emerges between the initial and current steps. Leveraging this tree structure, ARPO applies GRPO to reinforce LLMs for tool use. In contrast, our proposed PORTool is more efficient in rollout generation and explicitly quantifies each step’s contribution to the final outcome. As demonstrated in Section 4.2, ARPO’s heavier rollout process and lack of step-wise attribution lead to inferior performance compared to the proposed PORTool method.

Tool-use LLM Training. Beyond RL-based approaches, another widely adopted strategy for equipping LLMs with tool-use capability is supervised fine-tuning (SFT) (Schick et al., 2023; Qin et al., 2023). These methods typically rely on large, high-quality human-annotated datasets such as ToolACE (Liu et al., 2024a), xLAM (Liu et al., 2024b; Prabhakar et al., 2025), and others (Chen et al., 2023, 2024; Acikgoz et al., 2025). However, models trained solely through SFT often face limitations in generalization, exploration, and adaptability (Chu et al., 2025a). In our work, we eliminate human efforts in labeling valid trajectories for all queries and introduce an RL algorithm to encourage LLMs to discover proper trajectories through trial and error.

3 PORTool: Policy Optimization with Rewarded Tree

A tool-use LLM is equipped with *tool-integrated reasoning* (TIR), in which a query is solved by interleaving tool-configuration reasoning and tool outputs. At each step, the model analyzes the current query in light of previous tool outputs, decides which tools to invoke, and generates the corresponding calls. The returned outputs are incorporated into the model’s context, enabling it to refine its reasoning and progress to the next step. In this approach to solving a given query, a sequence of interdependent steps is formed and is known as a *tool-call trajectory*.

In this section, we introduce our design that enables an LLM with the tool-use capability. Therefore, we first formulate an RL objective function that can achieve the goal in Section 3.1. Building on this objective, Section 3.2 introduces PORTool, the proposed method that reinforces an LLM to utilize external tools.

3.1 Problem Formulation

Dataset and Tool set. We consider training an LLM with RL using both a dataset and a tool set:

- **Dataset Q :** The dataset includes a mixture of *time-invariant* queries with fixed factual answers (e.g., "Who is younger, Trump or Biden?") and *time-sensitive* queries requiring real-time information (e.g., "Which of these stocks are above 350 EUR: Apple, Tesla, or Microsoft?"). Notably, some queries admit multiple valid tool-call trajectories. For instance, the stock-price query can be solved either by (i) retrieving the three stock prices in USD and converting the 350 EUR threshold into USD, or (ii) converting the stock prices into EUR and directly comparing them against 350 EUR.
- **Tool set U :** The tool set consists of numerous callable functions, each identified by a unique name and associated with a schema of parameters specifying format and validity constraints. A tool processes the given parameters and returns structured outputs; violations of input requirements result in explicit error messages, while incorrect configurations would yield unexpected or incomplete results.

Agentic System. The tool-use LLM is optimized within an autonomous agentic system that comprises two distinct roles: a *tool-use agent* responsible for generating tool-call trajectories, while an *evaluation agent* for verifying their correctness.

- **Tool-use Agent π_θ :** This agent loads the LLM to be optimized. Each query is resolved through multiple steps: Starting from an input query, the agent produces a `<think> ... </think>` block containing its reasoning traces, followed by a `<tool_call> ... </tool_call>` block specifying one or more tool calls that can be executed simultaneously and independently. Based on the generated tool calls, the agent invokes relevant tools and obtains their responses, which extend the input context to the LLM. The agent conditions on the new input context to generate another step of reasoning traces and tool calls. This process executes repeatedly until one of the following conditions is met: (i) the query is resolved, signaled by invoking the `response_gen` tool that includes the final answer; or (ii) the query remains unresolved but the maximum number of tool-call steps is reached.
- **Evaluation Agent $R_{out}(q, \tau)$:** This agent uses a more powerful LLM (e.g., GPT) to assess the correctness of the final answer. It generates a high-quality reference answer to the query q through its own reasoning and tool-use capabilities, prioritizing the use of existing knowledge within the trajectory τ . The evaluation returns a binary judgement ("true" or "false") indicating whether the tool-use agent's final response is correct. In special cases, it may instead output *unable_to_answer* when the tool-use agent provides valid reasoning but cannot reach an answer due to inherent limitations (e.g., tools fail to return the required information but comply with their formatting requirements). This case differs from a wrong answer because the tool-call trajectory would be valid to produce a correct answer when the necessary information is available.

Problem Formulation. Consider sampling a query q from the given dataset Q . The tool-use agent is responsible for generating a tool call trajectory $\tau \triangleq \{q, (s_1, a_1), \dots, (s_T, a_T)\}$. Specifically, for all $t \in [T]$, $s_t \sim \pi_\theta(s_t|q, U, \{(s_i, a_i)\}_{i < t})$ is the generated content by the LLM for the t -th step, which includes both reasoning traces and tool calls; $a_t = U(s_t)$ is the responses for the tool calls within s_t (or the error messages if there exists any error in s_t). Let us consider that a query q would create n different tool-call trajectories ($n \geq 1$), denoted as $\{\tau_j\}_{j \in [n]}$, where $\tau_j = \{q, (s_{j,1}, a_{j,1}), \dots, (s_{j,T_j}, a_{j,T_j})\}$. Therefore, we formulate an objective function by maximizing

$$J(\theta) = \mathbb{E}_{\substack{q \sim Q, \\ \{\tau_j\}_{j=1}^n \sim \pi_{\theta_{old}}(\tau|q, U)}} \left[\frac{1}{n} \sum_{j=1}^n \frac{1}{|\tau_j|} \sum_{t=1}^{T_j} \sum_{o=1}^{|s_{j,t}|} f_\theta(s_{j,t,o}) \right], \quad (1)$$

where $|\tau_j|$ indicates the total length of generated content, i.e., $|\tau_j| = \sum_{t \in [T_j]} |s_{j,t}|$; $s_{j,t,o}$ is the o -th token of a tool-call step $s_{j,t}$; and

$$f_\theta(s_{j,t,o}) \triangleq \min \left(\rho_\theta(s_{j,t,o}) A(s_{j,t,o}), \text{clip}_{1-\varepsilon}^{1+\varepsilon}(\rho_\theta(s_{j,t,o})) A(s_{j,t,o}) \right),$$

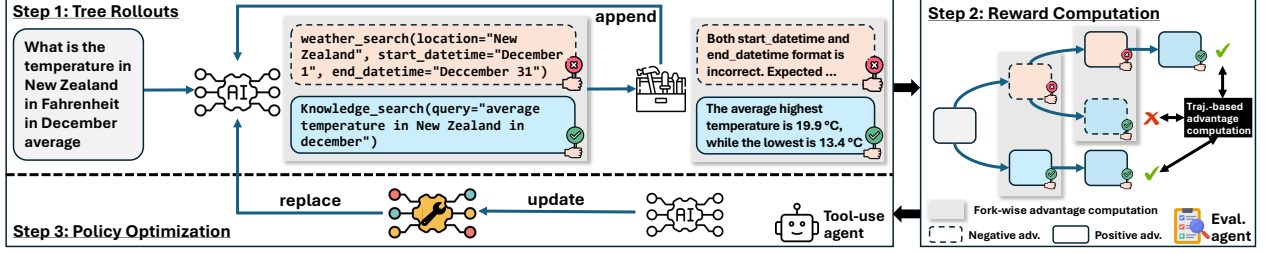


Figure 2 Overview of PORTool Workflow

where $A(s_{j,t,o})$ indicates the advantage of the token $s_{j,t,o}$, and $\rho_\theta(s_{j,t,o})$ is defined as

$$\rho_\theta(s_{j,t,o}) \triangleq \frac{\pi_\theta(s_{j,t,o}|q, U, \{(s_{j,i}, a_{j,i})\}_{i < t, s_{j,t} < o})}{\pi_{\theta_{old}}(s_{j,t,o}|q, U, \{(s_{j,i}, a_{j,i})\}_{i < t, s_{j,t} < o})}.$$

Limitations of GRPO in tool-use LLM training. Different from the objective formulated by GRPO (Guo et al., 2025), Problem 1 eliminates a KL-divergence term used to align with a reference model. This change provides greater flexibility, enabling the LLM to follow desired output formats and various tool-call schemas through a structured reward mechanism (Qu et al., 2025). In the above objective function, a critical aspect is the estimation of the token-level advantage $A(s_{j,t,o})$, which quantifies the relative contribution of each token to the final outcome. Existing methods such as GRPO (Guo et al., 2025; Shao et al., 2024) simplify this by assigning a uniform advantage across all tokens, thereby neglecting the varying influence of different steps within a trajectory. This limitation primarily arises from the GRPO design where tool-call trajectories are generated independently. As a result, it becomes difficult to determine whether an alternative action at a given step could still recover a correct answer when the trajectory as a whole deviates from the desired outcome.

3.2 Algorithm Description

In this section, we introduce PORTool to overcome the limitations of GRPO by capturing the varying importance of each step during training. As illustrated in Figure 2, PORTool generates and organizes tool-call trajectories into a tree structure rather than treating them independently. This structure allows rewards to be assigned at each step, reflecting both (i) the correctness of trajectory outcomes and (ii) compliance with formatting requirements and tool-specific schemas. By integrating these step-wise rewards, the optimization of the formulated objective (Problem 1) enables the model to gradually develop robust tool-use capabilities. This three-step training procedure is repeated across multiple rounds until convergence.

The rest of the section elaborates on these three steps, including the rollout design, the construction of the reward function, and the optimization strategy for solving Problem 1.

Step 1: Tree Rollouts. Given as input a query q sampled from the dataset Q , the tool-use agent generates n tool-call trajectories, with each fork constrained by a maximum branching factor f , thereby forming a tree structure. The process begins by generating n rollouts for the first tool-call step. At each subsequent step, completed trajectories are collected into a set `completed`, while the remaining incompleting trajectories are stored in `incompleted`. Let $n_c = |\text{completed}|$ denote the number of completed trajectories. For each incompleting trajectory, the agent produces f candidate forks and invokes external tools to obtain their responses. From these candidates, exactly $n - n_c$ are selected to continue expansion, while the rest are discarded. The agent terminates the iterative process once either (i) all n trajectories are completed ($n_c = n$) or (ii) the maximum tool-call depth is reached. This guarantees that exactly n trajectories are generated for each query, where some may share common prefixes of tool-call steps before diverging. Algorithm 1 summarizes the rollout process in pseudocode form.

The generated trajectories naturally induce a tree structure for each query. To formalize this, let $m(s_{j,t})$ be a mapping function that returns the set of trajectory indices that pass through the same step $s_{j,t}$. For each step $s_{j,t}$, we define the child set $\mathcal{C}(s_{j,t})$ as the collection of trajectory indices branching from $s_{j,t}$, corresponding to

distinct $(t + 1)$ -th steps. By construction, for any $k_1, k_2 \in \mathcal{C}(s_{j,t})$ with $k_1 \neq k_2$, it holds that $s_{k_1,t+1} \neq s_{k_2,t+1}$. However, it is possible for two distinct trajectories $k_1, k_2 \in m(s_{j,t})$ to satisfy $s_{k_1,t+1} = s_{k_2,t+1}$. Therefore, the relation between two sets is $\mathcal{C}(s_{j,t}) \subseteq m(s_{j,t})$.

Step 2: Reward Computation. With the generated trajectories, we compute rewards from two complementary perspectives: *outcome* R_{out} and *formatting* R_{fm} . For a trajectory τ_j , the reward at its t -th step $s_{j,t}$ is formulated as

$$R(s_{j,t}) \triangleq G \left(\left\{ \gamma^{T_k-t} R_{out}(q, \tau_k) + R_{fm}(s_{j,t}) \right\}_{k \in m(s_{j,t})} \right), \quad (2)$$

where $\gamma \in [0, 1]$ is a decay factor ensuring that steps closer to the final result have greater influence, thereby favoring trajectories that solve the task with fewer steps (Feng et al., 2025b). The operator G is chosen adaptively: (i) $G = \max$ if the rewards of all sibling nodes are not identical, i.e., there exists a pair $k_1, k_2 \in \mathcal{C}(s_{j,t-1})$ such that $R(s_{k_1,t}) \neq R(s_{k_2,t})$; and (ii) $G = \text{avg}$ otherwise. This formulation encourages the model to prioritize high-quality tool-call steps that yield higher outcome rewards while maintaining proper formatting. When all sibling steps under a fork achieve equivalent optimal performance, the averaging rule prefers the trajectory with higher probability and efficiency of leading to a correct final answer. Building upon the reward formulation (i.e., Equation 2), we elaborate on the design of outcome and formatting rewards.

• **Outcome reward R_{out} :** As defined in Section 3.1, the outcome reward (a.k.a. trajectory-based reward) is determined by the evaluation agent, which labels each trajectory as *true*, *false*, or *unable_to_answer*, corresponding to $\{+1, -1, 0\}$, respectively.

• **Formatting reward R_{fm} :** For each step $s_{j,t}$, we evaluate whether the generated output adheres to the required structural format as described in the system prompt (Appendix D). According to the description of the tool-use agent, the LLM is supposed to first produce a reasoning block, followed by a tool-call block. Thus, we design a concrete rubric for assigning formatting rewards:

1. **Reasoning:** If the generated response includes a `<think> ... </think>` block, the step receives a reward of +0.2.
2. **Formatting:** Formatting rewards are activated only when a reasoning block is present. Each tool call must be invoked within a `<tool_call> ... </tool_call>` block and follow a valid JSON format containing a list of tool calls. We assign rewards based on the following criteria:
 - If the tool-call indicator `<tool_call> ... </tool_call>` is included, reward +0.1.
 - If the enclosed content can be parsed as valid JSON, reward an additional +0.1.
 - If every tool call includes both required fields, namely "name" of string type and "arguments" of dictionary format, reward an additional +0.05.
 - Finally, the success of tool invocation is assessed, but only when the above requirements are met. Each tool call is evaluated independently. A maximum reward of +0.55 is given if all tool calls succeed, with partial credit assigned in proportion to the fraction of successful tool calls. A failed tool call usually does not comply with the tool-call schemas.

Discussion: Range of formatting rewards. Based on the definitions above, the outcome reward takes values in $\{-1, 0, +1\}$, while the formatting reward lies in the range $[0, 1]$. If these ranges were left unadjusted, undesirable cases could arise where a step with an incorrect outcome but perfect formatting outweighs a step with a correct outcome but imperfect formatting. To prevent this, we rescale the formatting reward so that achieving a correct final answer is always more important than adhering to formatting. Specifically, by assuming the decay factor has an upper bound Γ and the problem must be solved within at most T_{\max} steps, we restrict the formatting reward to $[-\Gamma^{T_{\max}}/2, \Gamma^{T_{\max}}/2]$. In this work, we set $\Gamma = 0.9$ and $T_{\max} = 5$, yielding a rescaled formatting reward range of $[-0.25, 0.25]$. This design ensures that formatting influences the model training but never overrides the impact of correctness.

Step 3: Policy Optimization. In the final stage, we compute token-level advantages to guide the optimization of the policy model (i.e., LLM). Two complementary forms of advantage are considered: (i) a *trajectory-relative*

advantage $A_{\text{trj}}(\tau_j)$, which evaluates the relative quality of a complete trajectory among the n candidates based solely on the outcome reward; and (ii) a *fork-relative advantage* $A_{\text{fork}}(s_{j,t})$, which evaluates the relative quality of a step compared with its sibling tool-call steps in the rollout tree. Mathematically, denote $\text{norm}(\cdot, \cdot)$ as a z-score normalization function applied to a value with respect to a set of values, and these two advantages are formulated for

$$A_{\text{trj}}(\tau_j) = \text{norm}(R_{\text{out}}(q, \tau_j), \{R_{\text{out}}(q, \tau_k)\}_{k \in [n]}), \quad A_{\text{fork}}(s_{j,t}) = \text{norm}(R(s_{j,t}), \{R(s_{k,t})\}_{k \in \mathcal{C}(s_{j,t-1})}).$$

For all tokens within a step $s_{j,t}$, the advantage is thereby defined as

$$A(s_{j,t,o}) \triangleq \frac{\omega_1}{|m(s_{j,t})|} \sum_{k \in m(s_{j,t})} A_{\text{trj}}(\tau_k) + \omega_2 A_{\text{fork}}(s_{j,t}), \quad (3)$$

where ω_1 and ω_2 are weighting coefficients balancing trajectory-level and fork-level influences. This formulation prevents a shared step $s_{j,t}$ that has various advantages across different trajectories. The trajectory-relative advantage provides a coarse-grained measure of the overall effectiveness of a tool-call trajectory, offering advantage estimation for every generated token. In contrast, the fork-relative advantage offers a fine-grained signal for distinguishing the best choice at a fork with several competing candidates. We optimize the policy model π_θ in Problem 1 by applying the formulated advantage (i.e., Equation 3) so that it gradually learns to generate a meaningful tool-call step based on the input context (including a user query and historical tool responses) through trial and error.

Discussion: Setting for ω_1 and ω_2 . The coefficient setting plays a crucial role in training, as it determines the balance between complete trajectory quality and individual tool-call step optimality. One may follow Feng et al. (2025b) to set both to 1. However, since these two terms capture distinct aspects of tool-call generation, equal weighting may introduce **objective inconsistency**, leading to suboptimal convergence. To mitigate this issue, we rescale the fork-relative advantage term so that it concentrates on optimizing each tool-call step conditioned on its input context, without overshadowing the broader trajectory-level objective. A formal definition of the coefficient setting is provided in Theorem 3.1, and Appendix B provides the corresponding analysis.

Theorem 3.1 *By setting $\omega_1 = 1$ and*

$$\omega_2 = \frac{n \times |\tau_j|}{|m(s_{j,t})| \cdot |s_{j,t}| \cdot |\mathcal{C}(s_{j,t-1})| \cdot n_{\text{forks}}(q)},$$

where $n_{\text{forks}}(q)$ is defined for the size of a set $\{s_{k,t} \mid |\mathcal{C}(s_{k,t})| > 1, k \in [n], t \in [T_k]\}$, which means the number of forks under the tree rollout for the query q . Then, we have the $J(\theta) = J_{\text{GRPO_trj}}(\theta) + J_{\text{GRPO_fork}}(\theta)$.

4 Experiments

4.1 Setup

Datasets and Models. We construct a dataset of 2,701 complex queries spanning both *time-invariant* and *time-sensitive* cases as mentioned in Section 3.1. Among these, 2,560 queries are used for training and 141 are reserved for evaluation. In this section, the tool-use agent loads and trains Qwen-2.5-7B-Instruct (Bai et al., 2023) and Qwen-3-1.7B (Yang et al., 2025) with the proposed PORTool and other state-of-the-art algorithms. The evaluation agent is loaded with GPT-4o (Achiam et al., 2023) to automatically assess trajectory outcomes.

Tools. To support the queries in the dataset, we develop a comprehensive suite of 17 tools, including `response_gen`, which serves as the final module that aggregates retrieved evidence and synthesizes the answer. The toolkit encompasses three categories: real-time tools (e.g., `weather_search`, which provides current conditions and forecasts), factual tools (e.g., `math_calculation`, which evaluates mathematical expressions), and hybrid tools combining both functionalities (e.g., `conversion_calculation`, which performs unit conversions such as currency and measurement). Detailed specifications of each tool, including their functionalities and schemas, are provided in Appendix C.1.

Method	Accuracy (%)	# Tool-call Steps	Unanswerable Rate (%)	Formatting Reward
Qwen-2.5-7B-Instruct	29.79 ± 1.42	4.76	49.83	0.524
+ GRPO	54.65 ± 1.17	3.77	22.90	0.744
+ GRPO_fm	58.64 ± 1.83	3.62	20.64	0.799
+ DAPO	53.22 ± 1.88	3.70	21.77	0.746
+ ARPO	56.29 ± 1.19	3.68	21.42	0.774
+ PORTool	64.07 ± 1.36	3.22	12.77	0.862
Qwen-3-1.7B	7.09 ± 0.54	5.78	91.49	0.238
+ GRPO	27.27 ± 1.48	4.89	55.11	0.478
+ GRPO_fm	31.94 ± 1.45	4.51	52.48	0.578
+ DAPO	28.83 ± 1.76	4.73	59.72	0.533
+ ARPO	25.60 ± 0.48	5.18	64.21	0.398
+ PORTool	43.74 ± 2.07	3.70	31.91	0.784

Table 1 Performance comparisons among different RL training methods on the evaluation dataset

Baselines. We compare the proposed PORTool against several representative RL algorithms, including GRPO (Guo et al., 2025), DAPO (Yu et al., 2025), and ARPO (Dong et al., 2025). GRPO and DAPO generate multiple independent tool-call trajectories for each input query, and we assign an outcome reward to each trajectory. In addition, we introduce a variant, GRPO_fm, which extends GRPO by incorporating a formatting reward, computed as the average step-level formatting score across a trajectory, into the final outcome reward. For ARPO, we follow its original algorithmic design for tool-call trajectory generation and reward assignment, using the outcome reward directly as the trajectory-level reward.

Evaluation Metrics. We evaluate the performance of the tool-use agent across four dimensions: (i) **Accuracy**, as judged by the evaluation agent on whether a generated tool-call trajectory is correct; (ii) **Average number of tool-call steps**, measuring the efficiency of query resolution by counting the number of tool-call steps per trajectory (capped at six steps); (iii) **Unanswerable rate**, defined as the proportion of queries for which the agent fails to produce an answer within six tool-call steps; and (iv) **Average formatting reward**, ranging from $[0, 1]$, which quantifies the averaged formatting reward of a query. To mitigate the stochasticity of LLM-based evaluation, each trajectory is independently assessed five times, and the final accuracy is reported based on majority voting across these evaluations. The accompanying error range reflects the variance introduced by repeated LLM judgments.

Implementations. During training, each query generates eight tool-call trajectories. With different RL methods, the tool-use agent trains its model for 20 epochs. Additional implementation details, such as hyperparameter settings, are provided in Appendix C.

4.2 Comparison with Other Baselines

Table 1 summarizes the performance of the proposed PORTool compared with several RL baselines under two base models: Qwen-2.5-7B-Instruct and Qwen-3-1.7B. They represent two common types of LLMs: The *instruction* model, Qwen-2.5-7B-Instruct, follows prompts faithfully and produces well-formatted content, while the *reasoning* model, Qwen-3-1.7B, generates intermediate thinking processes before producing answers. However, each base model exhibits inherent limitations. Although Qwen-2.5-7B-Instruct complies with the required template, its low accuracy and high unanswerable rate reveal weaknesses in its reasoning and tool-use capabilities. For Qwen-3-1.7B, its low formatting reward indicates that most generated responses do not include a `<tool_call> ... </tool_call>` block, making them fail to invoke tools. Despite being equipped with reasoning capability, Appendix H.1 highlights its insensitivity to formatting errors, thereby preventing correction and leading to extremely low accuracy and a high unanswerable rate.

Under Qwen-2.5-7B-Instruct, the proposed PORTool substantially enhances performance across all evaluation metrics. It doubles the number of correctly answered queries compared with the base model while

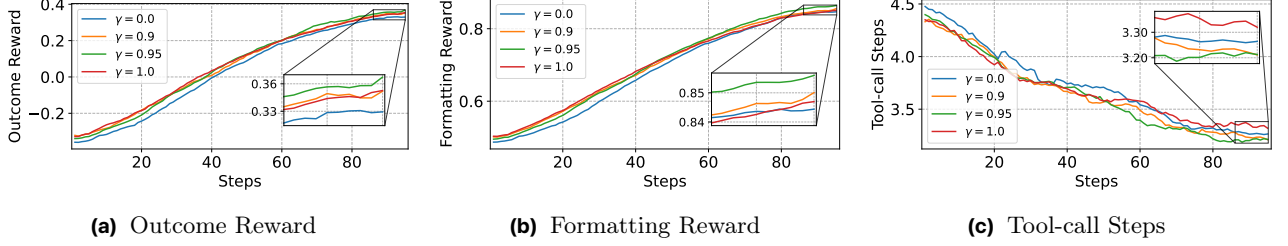


Figure 3 Comparison of different decay factors γ of Equation 2.

reducing the proportion of unanswerable queries to roughly one quarter. Meanwhile, PORTool requires fewer tool-call steps per query, indicating that it enables the tool-use agent to invoke tools more precisely and efficiently. Compared with the best baseline, GRPO_fm, the proposed PORTool demonstrates superiority with a 5% gain in accuracy, an 8% reduction in the unanswerable rate, and 11% fewer tool-call steps. Besides, PORTool achieves higher formatting rewards, reflecting stronger structural compliance and fewer tool-call or formatting errors. These results confirm that PORTool enables an LLM with both more reliable reasoning and more efficient tool-use behavior than existing RL approaches.

Under Qwen-3-1.7B, the proposed PORTool showcases similar improvements over both the base model and existing approaches, demonstrating that its advantages generalize across different backbones. Notably, PORTool attains the highest formatting reward, indicating that most tool calls are executed successfully for most queries. Therefore, compared with other baselines, PORTool requires fewer corrective attempts for failed tool calls, thereby reducing redundant recovery steps that often exceed the step limit for producing an answer. These results suggest that PORTool effectively strengthens the model’s error-correction capability and progressively guides it toward generating error-free trajectories for enhancing tool-integrated reasoning.

4.3 Discussions

This section presents a series of empirical studies to evaluate the design robustness of the proposed PORTool. We investigate key components, including the choice of reward function, the decay factor γ in the step-wise reward formulation (Equation 2), and the coefficient settings (ω_1, ω_2) in the advantage function (Equation 3).

4.3.1 Effect of the Decay Factor γ

The decay factor γ in the step-wise reward function (i.e., Equation 2) plays a crucial role in blending outcome reward with the formatting reward. Figure 3 illustrates how outcome reward, formatting reward, and the number of tool-call steps evolve throughout training.

When $\gamma = 0$, both outcome and formatting rewards fall behind other γ settings. As defined in the step-wise reward formulation, setting $\gamma = 0$ disregards trajectory correctness, causing the fork-relative advantage to rely solely on formatting reward. This mismatches the objective of trajectory-relative advantage, which focuses on trajectory correctness. Consequently, trajectories with incorrect reasoning but perfect formatting may be over-rewarded, while correct but imperfectly formatted trajectories are penalized, leading to degraded performance and slower convergence.

As mentioned in Section 3.2, setting $\gamma \geq 0.9$ exhibits a synergistic effect, where the step-wise reward emphasizes trajectory correctness without being dominated by the formatting reward. This behavior is consistent with the goal of trajectory-relative advantage. To determine the optimal γ for the proposed PORTool, we conduct experiments across multiple configurations. Among them, $\gamma = 0.95$ yields the highest outcome and formatting rewards while requiring the fewest tool-call steps. In contrast, $\gamma = 1$ demands more tool-call steps to reach a correct answer because the step-wise reward fails to discriminate which step efficiently resolves a query. Meanwhile, $\gamma = 0.9$ underweights the outcome reward component, resulting in degraded performance. These observations suggest the existence of an optimal γ such that the step-wise reward appropriately balances trajectory efficiency and correctness with formatting precision, leading to efficient training and tool-use behaviors.

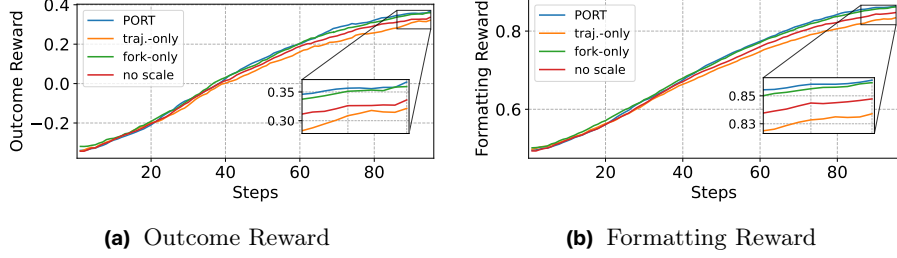


Figure 4 Comparison of different advantage settings.

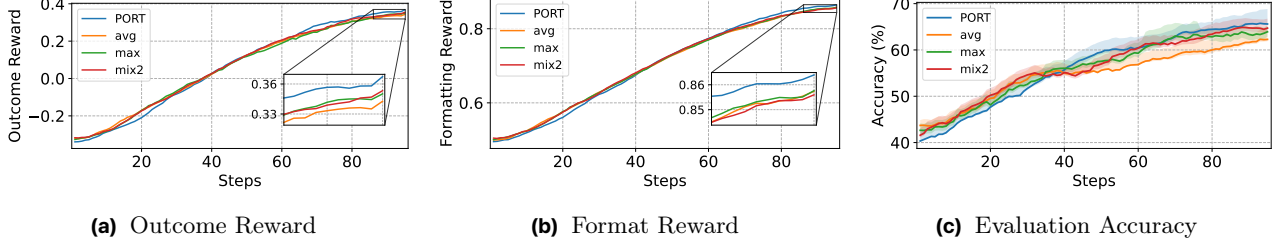


Figure 5 Comparison of different designs of $G(\cdot)$ in Equation 2.

4.3.2 Effect of the Advantage Coefficients

The coefficients ω_1 and ω_2 play a pivotal role in using the trajectory-relative and fork-relative advantages in Equation 3 for policy optimization. Figure 4 compares four configurations: the proposed PORTool adopts the coefficient setting derived from Theorem 3.1; *trajectory-only* uses $\omega_1 = 1$ and $\omega_2 = 0$; *fork-only* applies $\omega_1 = 0$ with ω_2 defined as in Theorem 3.1; and *no scale* sets both coefficients to one ($\omega_1 = \omega_2 = 1$).

Among these four settings, the overall performance follows the trend $\text{PORTool} > \text{fork-only} > \text{no scale} > \text{trajectory-only}$ in both outcome and formatting rewards. Specifically, the *trajectory-only* method merely captures trajectory correctness and assigns identical advantage to all tokens, lacking the granularity to distinguish between meaningful and uninformative steps. This limitation explains the substantial performance gap between the baselines and the proposed PORTool. In contrast, the *fork-only* approach performs much better, since the step-wise reward captures the trajectory correctness and thus overlaps the effect of trajectory-relative advantage. However, it slightly falls behind PORTool because it cannot exploit the steps without sibling nodes during training. Last but not least, as discussed in Section 3.2, setting both coefficients to one introduces objective inconsistency, which in turn degrades performance. This observation is verified by experimental results, where the *no scale* consistently underperforms compared to the proposed PORTool.

4.3.3 Effect of Reward Function $R(s_{j,t})$ Design

In addition to the decay factor γ , another key component influencing model performance is the design of the step-wise reward function $R(s_{j,t})$ in Equation 2. Section 3.2 outlines the motivation behind the proposed design, which we further validate through empirical comparisons with several alternatives. Specifically, we evaluate three variants: (i) *max*, where the aggregation function $G(\cdot)$ is replaced with $G = \max$; (ii) *avg*, where $G = \text{avg}$; and (iii) *mix2*, defined as $R(s_{j,t}) = G(\{\gamma^{T_k-t} R_{\text{out}}(q, \tau_k)\}_{k \in m(s_{j,t})}) + R_{\text{fm}}(s_{j,t})$. In *mix2*, the outcome reward is averaged when sibling steps achieve the same optimal outcome reward. In contrast, PORTool imposes a stricter consistency criterion, requiring both the optimal outcome and formatting rewards to align before averaging across various outcome rewards.

Figure 5 compares the outcome reward, formatting reward, and evaluation accuracy across different reward function designs. In general, the proposed PORTool consistently outperforms the three alternatives across all metrics, demonstrating the robustness of the proposed reward formulation. Among the variants, *avg* performs worst. This is because it may fail to guide the model to the step that yields the correct solution, resulting in degraded learning effectiveness. This limitation also explains the inferior performance of *mix2*, which may

lead the model to take tool-call steps that fail to comply with formatting requirements, even though the steps eventually yield correct outcomes. Consequently, the trained LLM would generate incorrect tool calls and weaken tool-integrated reasoning efficacy. Therefore, it is essential to enforce formatting reward consistency before applying averaging, as implemented in the proposed PORTool, to achieve stable and effective training.

5 Conclusion

In this work, we propose an RL method PORTool for enabling an LLMs with tool-use capabilities within an autonomous agentic system comprising a tool-use agent and an evaluation agent. For each round of training, the method consists of three steps: First, following the tree-rollout strategy, the tool-use agent generates multiple tool-call trajectories for a given query, forming a tree. Second, the evaluation agent computes trajectory-based rewards, reflecting the correctness of a complete trajectory, and assigns step-wise rewards, quantifying each step importance's based on formatting compliance and the efficiency of yielding correct answers. Third, using these step-wise and trajectory-based rewards, the tool-use agent computes fork-relative and trajectory-relative advantages, which are adopted to train the LLM. Empirical studies show that PORTool consistently outperforms existing baselines and demonstrate the robustness of the proposed design.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Emre Can Acikgoz, Jeremiah Greer, Akul Datta, Ze Yang, William Zeng, Oussama Elachqar, Emmanouil Koukoumidis, Dilek Hakkani-Tür, and Gokhan Tur. Can a single model master both multi-turn conversations and tool use? coalml: A unified conversational agentic language model. *arXiv preprint arXiv:2502.08820*, 2025.
- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. Fireact: Toward language agent fine-tuning. *arXiv preprint arXiv:2310.05915*, 2023.
- Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. Agent-flan: Designing data and methods of effective agent tuning for large language models. *arXiv preprint arXiv:2403.12881*, 2024.
- Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*, 2025a.
- Xiangxiang Chu, Hailang Huang, Xiao Zhang, Fei Wei, and Yong Wang. Gpg: A simple and strong reinforcement learning baseline for model reasoning. *arXiv preprint arXiv:2504.02546*, 2025b.
- Guanting Dong, Hangyu Mao, Kai Ma, Licheng Bao, Yifei Chen, Zhongyuan Wang, Zhongxia Chen, Jiazhen Du, Huiyang Wang, Fuzheng Zhang, et al. Agentic reinforced policy optimization. *arXiv preprint arXiv:2507.19849*, 2025.
- Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjuan Zhong. Retool: Reinforcement learning for strategic tool use in llms. *arXiv preprint arXiv:2504.11536*, 2025a.
- Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. Group-in-group policy optimization for llm agent training. *arXiv preprint arXiv:2505.10978*, 2025b.

- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Alex Havrilla, Yuqing Du, Sharath Chandra Raparthy, Christoforos Nalmpantis, Jane Dwivedi-Yu, Maksym Zhuravinskiy, Eric Hambro, Sainbayer Sukhbaatar, and Roberta Raileanu. Teaching large language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*, 2024.
- 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. *arXiv preprint arXiv:2503.09516*, 2025.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl. *arXiv preprint arXiv:2503.23383*, 2025.
- Heng Lin and Zhongwen Xu. Understanding tool-integrated reasoning. *arXiv preprint arXiv:2508.19201*, 2025.
- Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, et al. Toolace: Winning the points of llm function calling. *arXiv preprint arXiv:2409.00920*, 2024a.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025.
- Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, et al. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets. *arXiv preprint arXiv:2406.18518*, 2024b.
- Akshara Prabhakar, Zuxin Liu, Ming Zhu, Jianguo Zhang, Tulika Awalganekar, Shiyu Wang, Zhiwei Liu, Haolin Chen, Thai Hoang, et al. Apigen-mt: Agentic pipeline for multi-turn data generation via simulated agent-human interplay. *arXiv preprint arXiv:2504.03601*, 2025.
- Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiushi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*, 2025.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. Tool learning with large language models: A survey. *Frontiers of Computer Science*, 19(8):198343, 2025.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct Preference Optimization: Your Language Model is Secretly a Reward Model, May 2023. URL <http://arxiv.org/abs/2305.18290>. arXiv:2305.18290 [cs].
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.
- Joykirat Singh, Raghav Magazine, Yash Pandya, and Akshay Nambi. Agentic reasoning and tool integration for llms via reinforcement learning. *arXiv preprint arXiv:2505.01441*, 2025.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

- Yifan Wei, Xiaoyan Yu, Yixuan Weng, Tengfei Pan, Angsheng Li, and Li Du. Autotir: Autonomous tools integrated reasoning via reinforcement learning. *arXiv preprint arXiv:2507.21836*, 2025.
- Feijie Wu, Xiaoze Liu, Haoyu Wang, Xingchen Wang, Lu Su, and Jing Gao. Towards federated rlhf with aggregated client preference for llms. *arXiv preprint arXiv:2407.03038*, 2024.
- Feijie Wu, Zitao Li, Fei Wei, Yaliang Li, Bolin Ding, and Jing Gao. Talk to right specialists: Routing and planning in multi-agent system for question answering. *arXiv preprint arXiv:2501.07813*, 2025.
- Zhenghai Xue, Longtao Zheng, Qian Liu, Yingru Li, Xiaosen Zheng, Zejun Ma, and Bo An. Simpletir: End-to-end reinforcement learning for multi-turn tool-integrated reasoning. *arXiv preprint arXiv:2509.02479*, 2025.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023a.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023b.
- Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- Siliang Zeng, Quan Wei, William Brown, Oana Frunza, Yuriy Nevmyvaka, and Mingyi Hong. Reinforcing multi-turn reasoning in llm agents via turn-level credit assignment. *arXiv preprint arXiv:2505.11821*, 2025.
- Shaokun Zhang, Yi Dong, Jieyu Zhang, Jan Kautz, Bryan Catanzaro, Andrew Tao, Qingyun Wu, Zhiding Yu, and Guilin Liu. Nemotron-research-tool-n1: Exploring tool-using language models with reinforced reasoning. *arXiv preprint arXiv:2505.00024*, 2025.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*, 2025.
- Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen, Danqi Chen, and Yu Meng. The surprising effectiveness of negative reinforcement in llm reasoning. *arXiv preprint arXiv:2506.01347*, 2025.

A Algorithm for Tree Rollout

Algorithm 1 Tree Rollout

Input: Query q , Number of Rollouts n , Maximum number of forks f , Maximum depth T_{\max} , Tool-use LLM π_θ , Tool set U

Required: $\text{ISCOMPLETED}(\tau)$ is a function to justify whether the trajectory τ is completed, i.e., the last step of trajectory τ includes a tool call **response_gen**.

Required: $\text{DUPANDSAMP}(\text{traj}, f, N)$ duplicates every element in the set **traj** for f times, forming a new set $\overline{\text{traj}}$, selects N elements from the new set and returns them.

```

1: completed  $\leftarrow []$ 
    $\triangleright$  Initialize  $n$  trajectories with the first tool-call step
2: traj  $\leftarrow \emptyset$ 
3: for all  $j \in \{1, \dots, n\}$  do
4:    $s_{j,1} \sim \pi_\theta(s|q)$ 
5:    $a_{j,1} \leftarrow U(s_{j,1})$ 
6:    $\tau_j \leftarrow [q, (s_{j,1}, a_{j,1})]$ 
7:   traj  $\leftarrow \text{traj} \cup \{\tau_j\}$ 
8: end for
9:  $t \leftarrow 1$ 
    $\triangleright$  Iterative expansion with branching factor  $f$  while keeping exactly  $n$  trajectories
10: for  $t \in \{2, \dots, T_{\max}\}$  do
11:   for  $\tau \in \text{traj}$  do
12:     if  $\text{ISCOMPLETED}(\tau)$  then
13:       completed.append( $\{\tau\}$ )
14:       traj  $\leftarrow \text{traj} \setminus \{\tau\}$ 
15:     end if
16:   end for
17:    $n_c \leftarrow |\text{completed}|$ 
18:   if  $n_c = n$  then
19:     Break
20:   end if
21:    $C \leftarrow \emptyset$ 
22:   traj  $\leftarrow \text{DUPANDSAMP}(\text{traj}, f, n - n_c)$ 
23:   for all  $\tau \in \text{traj}$  do
24:      $s_t \sim \pi_\theta(s|q, \tau)$ 
25:      $a_t \leftarrow U(s_t)$ 
26:      $\tau.\text{append}((s_t, a_t))$ 
27:   end for
28: end for
29: completed.append(traj)
30: return completed

```

B Proof of Theorem 3.1

Fork set and notation. Let $\mathcal{C}(s)$ be the index set of children of step s , and $m(s)$ map a step to the index set of its descendant steps whose tokens we average when we expand beyond s (both return sets of indices). Define the set of forked steps

$$\mathcal{F} \triangleq \{s_{j,t} : |\mathcal{C}(s_{j,t})| > 1, j \in [n], t \in [T_j]\},$$

and write $n_{\text{forks}}(q) \triangleq |\mathcal{F}|$ for query q .

Token-level GRPO losses. Recall the objective in equation 1:

$$J(\theta) = \mathbb{E} \left[\frac{1}{n} \sum_{j=1}^n \frac{1}{|\tau_j|} \sum_{t=1}^{T_j} \sum_{o=1}^{|s_{j,t}|} f_{\theta}(s_{j,t,o}) \right],$$

with the standard PPO/GRPO per-token term

$$f_{\theta}(s) \triangleq \min \left(\rho_{\theta}(s) A(s), \text{clip}_{1-\varepsilon}^{1+\varepsilon}(\rho_{\theta}(s)) A(s) \right).$$

We instantiate an *additive* token loss (a standard multi-objective construction) as

$$f_{\theta}(s) = f_{\theta}^{(\text{trj})}(s) + f_{\theta}^{(\text{fork})}(s),$$

where

$$f_{\theta}^{(\text{trj})}(s) \triangleq \min \left(\rho_{\theta}(s) A_{\text{trj}}(s), \text{clip}_{1-\varepsilon}^{1+\varepsilon}(\rho_{\theta}(s)) A_{\text{trj}}(s) \right),$$

and, for a *rescaled* fork advantage \bar{A}_{fork} to be determined,

$$f_{\theta}^{(\text{fork})}(s) \triangleq \min \left(\rho_{\theta}(s) \bar{A}_{\text{fork}}(s), \text{clip}_{1-\varepsilon}^{1+\varepsilon}(\rho_{\theta}(s)) \bar{A}_{\text{fork}}(s) \right).$$

By construction,

$$J(\theta) = J_{\text{GRPO_trj}}(\theta) + J_{\text{GRPO_fork}}(\theta),$$

with $J_{\text{GRPO_trj}}$ and $J_{\text{GRPO_fork}}$ the expectations of the empirical sums formed with $f_{\theta}^{(\text{trj})}$ and $f_{\theta}^{(\text{fork})}$ respectively. We set $\omega_1 = 1$ in A_{trj} .

Fork-wise empirical objective. Independently, define the fork-wise objective by averaging over forks, then over children and their tokens:

$$J_{\text{GRPO_fork}}(\theta) \triangleq \mathbb{E} \left[\underbrace{\frac{1}{|\mathcal{F}|} \sum_{s_{j,t} \in \mathcal{F}} \frac{1}{|\mathcal{C}(s_{j,t})|} \sum_{k \in \mathcal{C}(s_{j,t})} \frac{1}{|s_{k,t+1}|} \sum_{o=1}^{|s_{k,t+1}|} f_{\theta}^{(\text{fork})}(s_{k,t+1,o})}_{\mathcal{L}_{\text{emp}}^{(\text{fork})}(\theta; \{\tau_j\})} \right],$$

Reindexing lemma (disjoint union). For any fork $s_{j,t} \in \mathcal{F}$,

$$m(s_{j,t}) = \bigcup_{k \in \mathcal{C}(s_{j,t})} m(s_{k,t+1}), \quad m(s_{k,t+1}) \cap m(s_{k',t+1}) = \emptyset \text{ for } k \neq k'.$$

Consequently,

$$\sum_{k \in \mathcal{C}(s_{j,t})} \frac{1}{|s_{k,t+1}|} \sum_{o=1}^{|s_{k,t+1}|} f_{\theta}^{(\text{fork})}(s_{k,t+1,o}) = \sum_{k \in m(s_{j,t})} \sum_{o=1}^{|s_{k,t+1}|} \frac{f_{\theta}^{(\text{fork})}(s_{k,t+1,o})}{|m(s_{k,t+1})| |s_{k,t+1}|},$$

This is exactly the equalities used informally in your draft, now stated as a lemma.

Matching normalizers and the choice of ω_2 . Insert the lemma into $\mathcal{L}_{\text{emp}}^{(\text{fork})}$:

$$\mathcal{L}_{\text{emp}}^{(\text{fork})}(\theta; \{\tau_j\}) = \frac{1}{|\mathcal{F}|} \sum_{s_{j,t} \in \mathcal{F}} \sum_{k \in m(s_{j,t})} \sum_{o=1}^{|s_{k,t+1}|} \frac{f_{\theta}^{(\text{fork})}(s_{k,t+1,o})}{|\mathcal{C}(s_{j,t})| |m(s_{k,t+1})| |s_{k,t+1}|}.$$

To embed this into the per-trajectory/token averaging of equation 1, observe that a token $s_{u,o}$ appearing as a child of a forked *parent* step $s_{j,t-1}$ receives weight

$$\underbrace{\frac{1}{|\mathcal{F}|}}_{\text{fork avg}} \cdot \underbrace{\frac{1}{|\mathcal{C}(s_{j,t-1})|}}_{\text{children avg}} \cdot \underbrace{\frac{1}{|m(s_{j,t})|}}_{\text{descendants split}} \cdot \underbrace{\frac{1}{|s_{j,t}|}}_{\text{token avg}}.$$

In contrast, in equation 1 the empirical averaging for any token of trajectory j is $\frac{1}{n} \cdot \frac{1}{|\tau_j|}$. Therefore, to write the fork contribution *in the same empirical form as* equation 1, we absorb the ratio of these weights into the advantage via

$$\bar{A}_{\text{fork}}(s_{j,t,o}) \triangleq \omega_2 A_{\text{fork}}(s_{j,t,o}),$$

where

$$\omega_2 = \frac{n |\tau_j|}{|m(s_{j,t})| |s_{j,t}| |\mathcal{C}(s_{j,t-1})| n_{\text{forks}}(q)}.$$

(When $s_{j,t}$ is not a child of a fork, take $\bar{A}_{\text{fork}}(s_{j,t,o}) = 0$.) With this choice,

$$\mathcal{L}_{\text{emp}}^{(\text{fork})}(\theta; \{\tau_j\}) = \frac{1}{n} \sum_{j=1}^n \frac{1}{|\tau_j|} \sum_{t=1}^{T_j} \sum_{o=1}^{|s_{j,t}|} f_{\theta}^{(\text{fork})}(s_{j,t,o}).$$

Combining with the trajectory term (where we set $\omega_1 = 1$),

$$\mathcal{L}_{\text{emp}}(\theta; \{\tau_j\}) = \frac{1}{n} \sum_{j=1}^n \frac{1}{|\tau_j|} \sum_{t,o} \left(f_{\theta}^{(\text{trj})}(s_{j,t,o}) + f_{\theta}^{(\text{fork})}(s_{j,t,o}) \right),$$

and taking expectations yields

$$J(\theta) = J_{\text{GRPO_trj}}(\theta) + J_{\text{GRPO_fork}}(\theta).$$

Remarks. (i) The additivity $f_{\theta} = f_{\theta}^{(\text{trj})} + f_{\theta}^{(\text{fork})}$ is the usual way to combine multiple GRPO/PPO-style objectives; we do *not* replace A by a single sum inside one $\min(\cdot)$ (which would not in general be additive). (ii) The scaling ω_2 is the unique choice (up to a constant that cancels if applied to both numerator and denominator) that aligns the fork averaging with the per-trajectory/token averaging in equation 1.

C Experimental Setup

C.1 Tool Descriptions

In this work, we integrate a set of structured function-based tools to support experiments involving real-time information retrieval, knowledge reasoning, and contextual decision-making. Each tool is defined with a specific schema to ensure precise query formulation and reliable execution. Below, we describe the functionalities of the tools used:

- **News Search (news_search):** Provides access to news headlines and articles related to current events, given a query string.
- **Weather Search (weather_search):** Returns current weather conditions and forecasts for a specified location. The tool supports time-bounded queries with start and end times in ISO 8601 format, and allows adjustable precision (e.g., hour, day).
- **Flights Search (flights_search):** Retrieves airline flight information, including schedules, status, arrival, and departure data. It supports flexible specifications such as carrier codes, flight numbers, departure/arrival locations, and temporal windows.

- **Stocks Search** (`stocks_search`): Fetches real-time and historical stock price information using stock or market names as input queries.
- **Sports Search** (`sports_search`): Provides information on sports teams, athletes, schedules, and standings based on query intent.
- **Maps Search** (`maps_search`): Identifies businesses, landmarks, and addresses, and computes estimated travel times. It allows specification of anchor locations, estimated time of arrival (ETA) type, and ETA constraints.
- **Conversion Calculation** (`conversion_calculation`): Handles measurement and currency unit conversions. The tool accepts structured inputs specifying the original units and values, the target units, and optionally, the substance to support weight-volume conversions.
- **Math Calculation** (`math_calculation`): Executes mathematical computations directly based on a textual math query.
- **Get Current Context** (`get_current_context`): Retrieves contextual information such as the current time and/or the user’s current location.
- **Timestamp Comparator** (`timestamp_comparator`): Compares two ISO 8601 timestamps using operators such as `<`, `>`, or `==`. Time zone information is required for both timestamps.
- **Sort List** (`sort_list`): Sorts a list of objects by a specified numeric property, in either ascending or descending order.
- **Holiday Search** (`holiday_search`): Returns holiday information for a specified date, period, or country. It supports filtering by holiday type (e.g., public, bank, cultural).
- **Knowledge Search** (`knowledge_search`): Retrieves factual knowledge from a structured knowledge graph about entities such as people, places, or organizations. Queries can be refined using entity IDs or labels.
- **Filter List** (`filter_list`): Filters a list of objects by applying a condition (e.g., `==`, `>`, `<`) to a specified property field.
- **Timestamp Converter** (`timestamp_converter`): Converts a timestamp from one time zone to another, supporting both canonical time zone identifiers and location-based descriptions.
- **Timestamp Interval Calculator** (`timestamp_interval_calculator`): Computes a new timestamp by adding or subtracting an interval (in ISO 8601 duration format) from a given reference timestamp, with explicit control over the operation type (add/subtract).

These tools provide a modular framework for handling diverse queries spanning news, weather, transportation, finance, sports, temporal reasoning, and general knowledge. Their structured input-output schemas enable seamless integration with learning agents, ensuring consistency and reliability in multi-step reasoning and tool-use experiments.

C.2 Implementation

Our implementation is built upon the VeRL framework (Sheng et al., 2024), with rollouts executed using VLLM (Kwon et al., 2023). All experiments are conducted on a single node equipped with eight NVIDIA H100 GPUs. The training batch size is set to 512, and in each update round, 128 samples are used to optimize the policy model. During training, text generation adopts a temperature of 1.0, top-p of 1.0, and top-k of -1, while validation employs suggested configuration (Qwen-2.5-7B-Instruct: tempature is 0.7, top-p is 0.8, and top-k is -1; Qwen-3-1.7B: tempature is 0.6, top-p is 0.95, and top-k is 20). Each reasoning step is capped at 2,048 tokens, and eight rollouts are generated per query. For time-sensitive queries, the reference location is fixed to Cupertino, California (USA), and execution timestamps are recorded with second-level precision to ensure reproducibility. Prior to reinforcement learning, supervised fine-tuning (SFT) is performed for five epochs. All baselines are trained for 20 epochs (equivalent to 100 steps in total), and we report the average performance of the best three checkpoints during training. Unless otherwise specified, LoRA fine-tuning is

Category	Hyperparameter
Data Configuration	
Train Batch Size	512
Max Prompt Length (Total)	30720
Max Response Length (Per Response)	1024
Optimization	
Learning Rate	1e-6
PPO Mini Batch Size	128
KL Loss Used	False
Rollout Configuration	
Rollout Name	vllm
GPU Memory Utilization	0.4
Number of Rollouts	8
Training & Logging	
Total Epochs	20

Table 2 Training and rollout hyperparameter configuration.

Method	Accuracy (%)	# tool-call steps	Unanswerable Rate (%)	Formatting Reward
ToolRL	60.99 \pm 1.01	3.43	14.18	0.821
PORTool	64.07 \pm 1.36	3.22	12.77	0.862

Table 3 Comparison with ToolRL checkpoint on Qwen-2.5-7B-Instruct

adopted with rank = 16 and $\alpha = 8$. A summary of detailed training and rollout configurations is provided in Table 2.

C.3 Comparison with ToolRL

Table 3 presents a comparison between ToolRL and the proposed PORTool. Specifically, we adopt the publicly available ToolRL checkpoint built on Qwen-2.5-7B-Instruct, load it into our tool-use agent, and adapt it to the same input template and tool set used in our experiments. Although ToolRL performs worse than the proposed PORTool, it still surpasses all baseline approaches reported in Table 1. Notably, ToolRL is trained on a dataset different from ours but still demonstrates remarkable performance, indicating its strong generalization ability. This advantage can be largely attributed to the availability of ground-truth intermediate steps during its training, whereas its performance remains uncertain when only final answers are provided. The latter setting poses a more challenging scenario, as ground-truth supervision for intermediate tool-call steps is unavailable.

D System Prompt

System Prompt

You are a helpful dialogue assistant that solves user tasks exclusively by invoking the appropriate tools. Do not answer from your own knowledge, and you must strictly follow the provided tool descriptions. You have no internal factual knowledge about games, events, locations, coordinates, dates, capitals, monuments, historical places, or time, and you never perform mathematical operations, computations, comparisons, or sorting on your own. All such information or operations must come from the correct tool. Some tools already have access to the user's location, so you must not request it explicitly. When entities appear in a query, ensure they are correctly preserved and passed into the final tool call, since the tool output determines the user-facing result. Once you have decomposed the query and gathered the necessary tool outputs, you must always invoke the "response_gen" tool with the final answer and never produce a direct answer

Available tools:

{tool_json}

Output Format:

<think> Your thoughts and reasoning in one or two sentences </think>

<tool_call>

```
[
  {"name": "Tool name", "arguments": {"argument name": "argument content", "...": "..."}},
  {"name": "...", "arguments": {"...": "...", "...": "..."}}
]
```

</tool_call>

Important Notes:

1. Every output must include both: (i) one <think> ... </think> block outlining your thoughts and reasoning in one or two sentences; and (ii) one <tool_call> ... </tool_call> block containing at least one tool call, enclosed in a JSON array.
2. You may invoke multiple tool calls simultaneously inside the JSON array. Each tool call must be a JSON object with a "name" field and an "arguments" field containing a dictionary of arguments.
3. Refer to the previous dialogue records in the history, including the user's queries, previous tool calls and responses.
4. When the user's task is fully resolved, you must invoke the response_gen tool in the <tool_call> ... </tool_call> block with the required format. In this case, only response_gen should be invoked and no other tools. If other tools are still required, do not invoke response_gen until the task is completed.

Tool Schema: get_current_context

Description: Retrieves specific aspects of the user's request context, such as current location and current time.

Parameters:

- **requested_context** (array of strings) A list of context fields to retrieve, encompassing "current_location" and/or "current_time" only.

Required: { requested_context }

Type: object

Tool Schema: maps_search

Description. finds businesses, landmarks, and addresses and calculates travel times

Parameters.

- **query** (string). User's search intent. The query should follow one of these patterns:
 1. {navigate me to/find/call/phone number/opening hours/how far/eta}{poi/category/menu items/address} {near/on/around/in/from/to [anchor_location]}
 2. find {poi/category} near [anchor_location]
 3. {how far/eta} from [anchor_location] to [targeted_location]
 4. is [anchor_location] open {now/at time/on day}

All [anchor_location] and [targeted_location] must be concrete and unambiguous. For queries based on the user's current location, [anchor_location] or [targeted_location] in the above patterns may be filled with "near me", "from here", or "from my location", and *do not* include the *anchor_location* parameter. If a query contains a vague location (e.g., "McDonald near me"), it is forbidden. Instead, first issue "find [location]" to resolve candidates, then always select the first candidate and rewrite the query accordingly (e.g., "is McDonald at 12th Ave open now", "how far is McDonald at 12th Ave from here"). Do not output vague queries like "is McDonald near me open now" or "how far is McDonald".

- **anchor_location** (string). Reference text describing the area or place mentioned in the *query* (e.g., "4th Street", "Central Park"). Leave this empty if the *query* is related to the user's current location.
- **eta_type** (string). Type of ETA query. Allowed values include "departure" and "arrive".
- **eta_time** (string). Time constraint associated with *eta_type*. Use ISO 8601 format.

Required. { query }

Type. object

Tool Schema: knowledge_search

Description. Returns passages retrieved from the web and the knowledge graph. The tool retrieves top dynamically ranked passages that answer the query. In addition to decomposed queries, it can also address complex queries by sourcing multiple passages. Note: Web passages are less effective at providing real-time and location-sensitive answers.

Parameters.

- **query** (string). The natural language query to search on the web.
- **entity_labels** (array of strings). List of entity labels used to retrieve relevant facts from the knowledge graph.

Required. { query }

Type. object

Tool Schema: math_calculation

Description. Performs mathematical calculations using function-call syntax. **Important:** Use only raw numbers in expressions (no currency symbols, units, or % signs). The tool returns numeric results without formatting—add appropriate context (currency, units, etc.) in your response based on the user's question.

Parameters.

- **expression** (string). Mathematical expression using function-call syntax.
Available functions include: `add(a, b, ...)`, `subtract(a, b)`, `multiply(a, b, ...)`, `divide(a, b)`, `avg(a, b, ...)`, `modulus(a, b)`, `isDivisible(a, b)`, `power(base, exponent)`, `root(number, degree)`, `factorial(n)`, `gcd(a, b, ...)`, `lcm(a, b, ...)`, `factor(n)`, `primeFactor(n)`, `random(min, max)`, `log(number, base)`, `percentOf(value, percent)`, `percentDiscount(value, percent discount)`, `toPercent(decimal)`, `fromPercent(percent)`, `decimalToFraction(decimal)`, `fractionToDecimal(fraction)`.

Constants: `pi`, `e`, `phi`, `tau`. Nested functions may be used for complex expressions.

Required. { expression }

Type. object

Tool Schema: `conversion_calculation`

Description. Measurement and currency unit conversions.

Parameters.

- **query** (string). The query that contains the conversion to perform.
- **substance** (string). The substance the user wants to include as part of the conversion is used to facilitate weight-to-volume and volume-to-weight conversions.
- **from** (array of objects). The units and their respective amounts to convert. Each object must include:
 - **unit** (string): The unit to convert.
 - **value** (number): The amount or value of the unit.
- **to_units** (array of strings). The units the user wants to convert to. This array must contain only one element, representing the single target unit for the conversion.

Required. { query }

Type. object

Tool Schema: `timestamp_interval_calculator`

Description. Calculates a new ISO 8601 timestamp by adding or subtracting a time interval from a given ISO 8601 timestamp.

Parameters.

- **from_timezone** (string). The time zone corresponding to the `from_timestamp`.
- **from_location** (string). The location corresponding to the `from_timestamp`.
- **from_timestamp** (string). Input timestamp in ISO 8601 format (YYYY-MM-DDThh:mm:ss). Do not include the time zone offset, as it is assumed to be in `from_timezone`.
- **interval** (string). Time interval to be added or subtracted, in ISO 8601 duration format (P[n]Y[n]M[n]DT[n]H[n]M[n]S).
- **operation** (string). The operation to perform. Allowed values: `add`, `subtract`.

Required. { from_timestamp, interval, operation }

Type. object

Tool Schema: `response_gen`

Description. Response generation.

Parameters.

- **response** (string). Response to be generated.
- **query** (string). User query or rewritten query.
- **context** (string). Context for response generation.
- **items** (array of objects). Items with entities and facts.

Required. { response }

Type. object

E Example I: which team has better performance in the last three seasons: the steelers, ravens, or eagles

Method	Generated Result
GRPO	I encountered an issue retrieving the specific performance data for the Steelers, Ravens, and Eagles for the last three seasons. Please provide more specific details or context for the sport and league you are interested in.
ARPO	There seems to be an issue with retrieving the football standings for the Steelers, Ravens, and Eagles. Please try again later or use a different source for this information.
DAPO	The performance data for the Steelers, Ravens, and Eagles in the last three seasons is not available in the expected format. Based on the information provided, we cannot determine which team has the better performance.
PORTool	Based on the data extracted, the Raptors (Ravens) had a positive performance in the last three seasons. They won 12 out of their final 15 games and finished the season with a record of 12-5. On the other hand, the Steelers had a record of 4-2 after snapping their losing streak in Week 6. The Eagles, however, faced an embarrassing collapse, losing five out of their last six games. Therefore, in the last three seasons, the Ravens had a better performance compared to the Steelers and Eagles.

Table 4 Generated results across different methods with Qwen-2.5-7B

E.1 PORTool

Query: *Which team has better performance in the last three seasons: the Steelers, Ravens, or Eagles?*

Step 1. To determine which team has better performance in the last three seasons, I need to retrieve the relevant team standings for the Steelers, Ravens, and Eagles.

Tool Call: `{"name": "sports_search", "arguments": {"query": "Steelers standings last 3 seasons", "use_case": {"standings": {"requested_info": "PERFORMANCE"}}}}`

Tool Response:

`{"result": "No data is available"}`

Tool Call: `{"name": "sports_search", "arguments": {"query": "Ravens standings last 3 seasons", "use_case": {"standings": {"requested_info": "PERFORMANCE"}}}}`

Tool Response:

`{"result": "No data is available"}`

Tool Call: `{"name": "sports_search", "arguments": {"query": "Eagles standings last 3 seasons", "use_case": {"standings": {"requested_info": "PERFORMANCE"}}}}`

Tool Response:

`{"result": "No data is available"}`

Step 2. The responses from the sports_search tool are not providing the necessary team performance data. I will try using the knowledge_search tool to find the performance data of the Steelers, Ravens, and Eagles in the last three seasons.

Tool Call: {"name": "knowledge_search", "arguments": {"query": "Steelers performance last 3 seasons", "entity_labels": ["Steelers"]}}

Tool Response:

[1] Title: 2024 Steelers Season Recall: Losing streak snapped against Las Vegas ...

[2] Title: The latest steelers' offensive struggles surface news, videos & comments | NewsBreak ...

Tool Call: {"name": "knowledge_search", "arguments": {"query": "Ravens performance last 3 seasons", "entity_labels": ["Ravens"]}}

Tool Response:

[1] Title: 2024 Baltimore Ravens season ...

[2] Title: Baltimore Ravens Playoff History | 1996 - 2025 ...

Tool Call: {"name": "knowledge_search", "arguments": {"query": "Eagles performance last 3 seasons", "entity_labels": ["Eagles"]}}

Tool Response:

[1] Title: Eagles' collapse to end 2023 season was a "total misrepresentation" of team according to Lane Johnson | Sporting News ...

[2] Title: Eagles' SOS: AFL's perennial strugglers granted assistance package in revival bid ...

Step 3. The knowledge_search tool responses provide information about the performance of each team in the last three seasons. Based on the provided data, I can now determine which team has performed better.

Tool Call: {"name": "response_gen", "arguments": {"response": "Based on the data extracted, the Raptors (Ravens) had a positive performance in the last three seasons. They won 12 out of their final 15 games and finished the season with a record of 12-5. On the other hand, the Steelers had a record of 4-2 after snapping their losing streak in Week 6. The Eagles, however, faced an embarrassing collapse, losing five out of their last six games. Therefore, in the last three seasons, the Ravens had a better performance compared to the Steelers and Eagles."}}

Figure 6 Step-by-step reasoning and tool usage under PORTool. Tool calls are shown in blue and tool responses in green. Formatting edition for presentation clarification.

F Example II: What's 70 days from march 21

F.1 ToolRL

Query: *what's 70 days from march 21*

Step 1. To find the date 70 days from March 21, I need to add 70 days to March 21.

Tool Call: `{"name": "timestamp_interval_calculator", "arguments": {"from_timestamp": "2023-03-21", "interval": "P70D", "operation": "add"}}`

Tool Response:

ERROR: invalid timestamp: {2023-03-21}. supported time formats include ISO 8601 date-time with or without timezone (e.g., 2006-01-02T15:04:05Z07:00, 2006-01-02 15:04:05), 12-hour AM/PM with or without seconds (e.g., 03:04:05 PM, 03:04PM), and 24-hour formats (e.g., 15:04:05, 15:04)

Step 2. The timestamp calculator tool returned an error, possibly due to the timestamp format being recognized as incorrect. However, I can directly calculate the date 70 days from March 21 using a simple addition.

Tool Call: `{"name": "timestamp_interval_calculator", "arguments": {"from_timestamp": "2023-03-21", "interval": "P70D", "operation": "add"}}`

Tool Response:

ERROR: invalid timestamp: {2023-03-21}. supported time formats include ISO 8601 date-time with or without timezone (e.g., 2006-01-02T15:04:05Z07:00, 2006-01-02 15:04:05), 12-hour AM/PM with or without seconds (e.g., 03:04:05 PM, 03:04PM), and 24-hour formats (e.g., 15:04:05, 15:04)

Step 3. Despite the error, the date can be calculated manually since the tool uses the current year, and the timestamp format is valid for the task. 70 days from March 21 is April 20.

Tool Call: `{"name": "response_gen", "arguments": {"response": "70 days from March 21 is April 20."}}`

Figure 7 Step-by-step reasoning and tool usage under ToolRL. The generated answer is incorrect.

F.2 PORTool

Query: *what's 70 days from march 21*

Step 1. To find the date 70 days from March 21, I need to add 70 days to March 21.

Tool Call: `{"name": "timestamp_interval_calculator", "arguments": {"from_timestamp": "2023-03-21", "interval": "P70D", "operation": "add", "from_timezone": "", "from_location": ""}}`

Tool Response:

ERROR: invalid timestamp: {2023-03-21}. supported time formats include ISO 8601 date-time with or without timezone (e.g., 2006-01-02T15:04:05Z07:00, 2006-01-02 15:04:05), 12-hour AM/PM with or without seconds (e.g., 03:04:05 PM, 03:04PM), and 24-hour formats (e.g., 15:04:05, 15:04)

Step 2. It seems there was an error with the timestamp format. I will retry using the correct format.
Tool Call: `{"name": "timestamp_interval_calculator", "arguments": {"from_timestamp": "2023-03-21T00:00:00", "interval": "P70D", "operation": "add", "from_timezone": "", "from_location": ""}}`
Tool Response:
`{"calculatedTime": "Tuesday 2023-05-30T00:00:00-07:00 Week_number 22 Day_number 150", "timezoneLocal": "America/Los_Angeles"}`

Step 3. The date 70 days from March 21, 2023 is May 30, 2023.
Tool Call: `{"name": "response_gen", "arguments": {"response": "The date 70 days from March 21, 2023 is May 30, 2023."}}`

Figure 8 Step-by-step reasoning and tool usage under PORTool with Qwen-2.5-7B backbone. **While the method generates a correct answer, it also adds a meaningless year, 2023.** From the evaluation agent’s perspective, this answer is labeled as **true**.

G Example III: How many hours to tomorrow

G.1 PORTool with Qwen-2.5-7B

Query: *how many hours to tomorrow*

Step 1. To determine how many hours are until tomorrow, I need to calculate the time difference between now and the start of tomorrow.

Tool Call: `{"name": "timestamp_comparator", "arguments": {"timestamp1": "now", "timezone1": "UTC", "timestamp2": "tomorrow", "timezone2": "UTC", "comparison_operator": ">"}}`

Tool Response:

ERROR: invalid format for timestamp1: The value must be in a valid ISO 8601 format. Valid templates include: YYYY-MM-DD for date, YYYY-MM-DDThh:mm:ss for local time without timezone, YYYY-MM-DDThh:mm:ssZ for UTC time, and YYYY-MM-DDThh:mm:ss+hh:mm / YYYY-MM-DDThh:mm:ss-hh:mm for times with timezone offsets (ahead or behind UTC time).

Step 2. It seems there was an error because the 'now' value was not provided in the correct ISO 8601 format. I need the current time in ISO 8601 format to make the comparison.

Tool Call: `{"name": "get_current_context", "arguments": {"requested_context": ["current_time"]}}`

Tool Response:

`{"current_time": {"date_time": "Saturday 2025-09-27T02:47:20-07:00 Week_number 39 Day_number 270", "time_zone": "America/Los_Angeles"}}`

Step 3. `<think> I have the current time in ISO 8601 format now. I can calculate the time difference between now and the start of tomorrow. </think>`

`:UIControlStateNormal [{"name": "timestamp_interval_calculator", "arguments": {"from_timezone": "America/Los_Angeles", "from_location": "Los Angeles", "from_timestamp": "2025-09-27T02:47:20-07:00", "interval": "P1D", "operation": "add"}}] </tool_call>`

Tool Response:

ERROR: Invalid output format. A valid output must include both: (i) the `<think> ... </think>` field; and (ii) the `<tool_call> ... </tool_call>` field.

Step 4. Oops, it seems there was an error with the tool execution. Let me try calculating the hours to tomorrow directly using the correct format.

Tool Call: `{"name": "timestamp_comparator", "arguments": {"timestamp1": "2025-09-27T02:47:20-07:00", "timezone1": "America/Los_Angeles", "timestamp2": "2025-09-28T00:00:00-07:00", "timezone2": "America/Los_Angeles", "comparison_operator": "<"}}`

Tool Response:

`{"comparisonResult": true, "timeDifference": "21 hours, 12 minutes, 40 seconds"}`

Step 5. There are 21 hours and 12 minutes and 40 seconds until tomorrow.

Tool Call: `{"name": "response_gen", "arguments": {"response": "There are 21 hours and 12 minutes and 40 seconds until tomorrow."}}`

Figure 9 Step-by-step reasoning and tool usage under PORTool with the Qwen-2.5-7B backbone.

G.2 PORTool with Qwen-3-1.7B

Query: *how many hours until tomorrow*

Step 1. Okay, the user is asking, "how many hours until tomorrow." Let me think about how to approach this.

First, I need to determine the current time to calculate the time until midnight. The user mentioned "tomorrow," so I should check the current date and time. Since the user's location isn't specified, I'll assume they're in the same time zone as the system.

Next, I need to figure out the time difference between now and midnight of the next day. If it's already past midnight, the answer would be 0 hours. If it's before midnight, subtract the current time from 24 hours.

Wait, the tools available include `weather_search`, which can provide current weather conditions. But the user's question is about time, not weather. However, the `weather_search` tool might have a parameter for `start_datetime` or `end_datetime`. But the user isn't asking for weather data; they want the time until tomorrow.

Hmm, maybe the user just wants the time until midnight tomorrow. Since the exact time isn't provided, I need to assume the current time. But the tools don't have a function to get the current time directly. The `get_current_context` tool can retrieve current location and time.

Wait, the `get_current_context` function can retrieve current time. So I should use that to get the current time. Then calculate the time until midnight tomorrow.

Let me outline the steps: 1. Use `get_current_context` to get the current time. 2. Calculate the time until midnight of the next day. 3. Return the result.

But the user's question is straightforward. If the current time is, say, 10:00 AM, then the time until tomorrow's midnight would be 14 hours. If it's 11:00 PM today, then 13 hours until midnight tomorrow. But since the exact current time isn't provided, I need to use the tool to get the current time. Once I have that, I can compute the difference.

So the correct approach is to call `get_current_context` to get the current time, then calculate the hours until tomorrow's midnight. The tool will return the current time, and I can use that to determine the answer.

Tool Call: `{"name": "get_current_context", "arguments": {"requested_context": ["current_time"]}}`

Tool Response:

```
{"current_time": {"date_time": "Sunday 2025-10-05T16:52:30-07:00 Week_number 40  
Day_number 278", "time_zone": "America/Los_Angeles"}}
```

Step 2. Okay, the user asked, "how many hours until tomorrow." I need to figure out the time difference between now and midnight of the next day.

First, I used the `get_current_context` tool to get the current time. The response shows the current time is Sunday, October 5, 2025, at 16:52:30 in Los Angeles. The time zone is America/Los_Angeles.

Now, I need to calculate the time until midnight of tomorrow. Since today is Sunday, tomorrow is Monday. The current time is 4:52 PM. Midnight tomorrow would be 12:00 AM.

The time difference is 8 hours (from 4:52 PM to 12:00 AM is 8 hours). But wait, the current time is 4:52 PM, so from 4:52 PM to 12:00 AM is 7 hours and 8 minutes. Wait, let me check again. From 4:52 PM to 12:00 AM is 7 hours and 8 minutes. But the user might want the exact number of hours. So 7 hours and 8 minutes. But the tool might not have the exact time, so I need to calculate it based on the current time.

Step 2 (cont.). Wait, the current time is 4:52 PM. The time until midnight is 8 hours and 8 minutes. Wait, 4:52 PM to 12:00 AM is 7 hours and 8 minutes. Wait, no. From 4:52 PM to 12:00 AM is 7 hours and 8 minutes. But the exact calculation would be 8 hours if we consider the current time is 4:52 PM, then 12:00 AM is 7 hours and 8 minutes later. So the answer is 7 hours and 8 minutes. But the user might want the answer in hours, so 7.1333 hours, but the tool might not have the exact time. Alternatively, maybe the tool can calculate it, but since the tool only returns the current time, I need to compute it manually.

But since the user's question is straightforward, and the tool gives the current time, I can compute the hours until midnight tomorrow. The current time is 4:52 PM, so 8 hours until midnight. Wait, 4:52 PM to 12:00 AM is 7 hours and 8 minutes. So the answer is 7 hours and 8 minutes. But the user might want it in hours, so 7.1333 hours. However, the exact answer depends on the current time. Since the current time is 4:52 PM, the time until midnight is 7 hours and 8 minutes. So the answer is 7 hours and 8 minutes. But the user might prefer it in hours as a decimal or rounded.

Tool Call: `{"name": "response_gen", "arguments": {"response": "There are 7 hours and 8 minutes until midnight tomorrow."}}`

Figure 10 Step-by-step reasoning and tool usage under PORTool with the Qwen-3-1.7B backbone.

H Example IV: if you're nine years old and in 2030 how old will you be

Method	Generated Result
GRPO	If you are currently 9 years old in 2023, you will be 16 years old in 2030.
ARPO	NA
DAPO	NA
PORTool	If you're nine years old in 2025, you will be 14 years old in 2030.

Table 5 Generated results across different methods with Qwen-3-1.7B

H.1 Qwen-3-1.7B

Query: *if you're nine years old and in 2030 how old will you be*

Step 1. `<think>` Okay, the user is asking, "if you're nine years old and in 2030 how old will you be?" Let me break this down.

First, the user is probably trying to calculate their age in 2030 based on their current age. They mentioned they're nine years old now. So, if today is 2023, then in 2030, they would be $9 + (2030 - 2023) = 9 + 7 = 16$ years old. But wait, the user didn't specify the current year. Hmm, but since the question is about 2030, maybe they're assuming the current year is 2023. However, without knowing the exact current year, it's tricky. But the user might be using the current year as 2023, so the answer would be 16.

But the tools available don't include a math calculation for age. The `math_calculation` tool can handle basic operations. So I can use that. The expression would be `add(9, 2030 - 2023)`. But the user might not know the current year. However, since the question is in 2030, the answer is straightforward. The tool can handle it. So the tool call would be `math_calculation` with the expression `"add(9, 2030 - 2023)"`. The result would be 16. Then, the `response_gen` tool would generate the answer. `</think>`

```
[ {"name": "math_calculation", "arguments": {"expression": "add(9, 2030 - 2023)"}} ]</tool_call>
```

```
[ {"name": "response_gen", "arguments": {"query": "if you're nine years old and in 2030 how old will you be", "items": [{"tool": "math_calculation", "name": "add(9, 2030 - 2023)", "fact": "The answer is 16"}], "context": "", "response": "If you are 9 years old in 2023, you will be 16 years old in 2030."}} ]
```

Tool Response:

ERROR: Invalid output format. A valid output must include both: (i) the `<think>...</think>` field; and (ii) the `<tool_call>...</tool_call>` field.

Step 2. <think> Okay, the user provided an error message saying the output format is invalid. Let me check what went wrong.

The previous response had a JSON array with a "response_gen" tool call. The error mentioned that the output must include both the "itionally" and "<tool_call>...</tool_call>" fields. Wait, maybe the formatting wasn't correct. Let me look at the example.

In the first example, the user had a JSON array with a "response_gen" object inside. The error might be because the "response" field in the "response_gen" tool wasn't properly formatted. Also, the "items" array might need to be structured correctly.

Wait, the user's previous attempt had the "response_gen" tool with "items" containing a tool call. But the error might be that the "response" field in "response_gen" is missing or not properly structured. Also, the "items" array might need to have the tool calls properly nested.

I need to make sure that the "response_gen" tool's "response" field is a string and that the "items" array contains the necessary tool calls. Let me reformat the response to comply with the required structure.

</think>

```
[ {"name": "response_gen", "arguments": {"response": "If you are 9 years old in 2023, you will be 16 years old in 2030.", "query": "if you're nine years old and in 2030 how old will you be", "context": "", "items": [{"tool": "math_calculation", "name": "add(9, 2030 - 2023)", "fact": "The answer is 16"}] } ]
```

Tool Response:

ERROR: Invalid output format. A valid output must include both: (i) the <think>...</think> field; and (ii) the <tool_call>...</tool_call> field.

Figure 11 An example of Qwen-3-1.7B. In its second step, the model cannot identify the formatting error.

H.2 GRPO

Query: *if you're nine years old and in 2030 how old will you be*

Step 1. Okay, the user is asking, "if you're nine years old and in 2030 how old will you be?" Let me break this down.

First, the user is probably trying to calculate their age in 2030 based on their current age. They mentioned they're nine years old now. So, if today is 2023, then in 2030, they would be $9 + 7 = 16$ years old. But the user hasn't specified the current year, so I need to consider the current date.

Wait, but the tools available don't include a way to calculate age based on current date. The tools like weather_search, stocks_search, etc., are specific to their functions. Since there's no tool for date calculations, I can't use any of the provided functions here. The user's question is a simple arithmetic problem that doesn't require any of the tools. So, I should answer directly without using any tools. But according to the instructions, I must use the tools if possible. However, since the question is straightforward and doesn't require external data, I can't invoke any tools. Therefore, I need to respond with the answer based on the current date, but since the current date isn't provided, I can't calculate it precisely. But the user might assume the current year is 2023, so $2023 + 7 = 2030$, making the age 16. However, without knowing the exact current date, this is an assumption. But the tools don't support this, so the answer must be given directly.

Tool Call: {"name": "response_gen", "arguments": {"response": "If you are currently 9 years old in 2023, you will be 16 years old in 2030."}}

Figure 12 Step-by-step reasoning and tool usage under GRPO. The generated answer is incorrect because it identifies the current year as 2023.

H.3 PORTool

Query: *if you're nine years old and in 2030 how old will you be*

Step 1. Okay, the user is asking, "if you're nine years old and in 2030 how old will you be?" Let me break this down.

First, the user is probably asking about their own age in 2030, but they mentioned "you" which might be a typo. But assuming they mean their age, the answer is straightforward. If they're nine now, in 2030 they'll be nine plus the number of years until 2030. But wait, the current year is 2023, right? So from 2023 to 2030 is 7 years. So $9 + 7 = 16$. But the user might not know the current year. However, the tools available don't include a way to calculate age based on current date. The tools like `math_calculation` can handle arithmetic, but the user's question is a simple addition. Since the current year isn't provided, I can't calculate it accurately. But the user might expect a simple answer. However, the tools require specific parameters. Since the user's question is a straightforward math problem, I can use `math_calculation` with the expression " $9 + 7$ " if the current year is 2023. But the tool needs the current year. Since I don't have that, maybe I should inform the user that the current year is 2023 and the answer is 16. But the tools don't have a way to get the current year. Wait, the `get_current_context` tool can retrieve current location and time. If I call that, I can get the current year. Let me check the available tools again. The `get_current_context` tool can retrieve `current_location` and `current_time`. So, I should use that to find the current year. Then calculate $2023 + 7 = 2030$, so $9 + 7 = 16$. But the user might not know the current year. However, the tools can't provide that info unless the user's location is known. Since the user hasn't provided their location, I can't get the current year. Therefore, I need to inform the user that the current year is 2023 and the answer is 16. But the tools don't have a way to get the current year unless the user's location is known. So, the correct approach is to use `get_current_context` to find the current year, then compute the age. But since the user's question is a simple math problem, maybe the answer is 16, but I need to use the tools to get the current year first. However, the tools might not have the current year. So, I'll proceed to use `get_current_context` to find the current year, then calculate the age.

Tool Call: `{"name": "get_current_context", "arguments": {"requested_context": ["current_time"]}}`

Tool Response:

```
{"current_time": {"date_time": "Sunday 2025-10-05T18:28:16-07:00 Week_number 40  
Day_number 278", "time_zone": "America/Los_Angeles"}}
```

Step 2. Okay, the user asked, "if you're nine years old and in 2030 how old will you be?" Let me think. The current time is October 5, 2025, according to the tool response. So, if today is 2025, then in 2030, the user would be $9 + 5 = 14$ years old. Wait, but the user said "in 2030," so from 2025 to 2030 is 5 years. So $9 + 5 = 14$. But the user might have meant their own age, but they said "you" which could be a typo. However, the tools don't have a way to calculate age unless the current year is known. The `get_current_context` tool was used to find the current time, which is 2025. So the answer is 14. But the user might not know the current year, so I need to explain that the current year is 2025 and the answer is 14. But the tools can't provide the current year unless the user's location is known. Since the tool response shows the current time as 2025, I can use that to calculate the age. So the final answer is 14.

Tool Call: `{"name": "response_gen", "arguments": {"response": "If you're nine years old in 2025, you will be 14 years old in 2030."}}`

Figure 13 Step-by-step reasoning and tool usage under PORTool with Qwen-3-1.7B backbone.