

AT²PO: Agentic Turn-based Policy Optimization via Tree Search

Zefang Zong^{1*}, Dingwei Chen^{1,2*}, Yang Li¹, Qi Yi¹, Bo Zhou¹,
Chengming Li³, Bo Qian¹, Peng Chen¹, Jie Jiang^{1†}

¹Tencent Inc ²Sun Yat-Sen University ³Shenzhen MSU-BIT University
{willzong, thomasyngli}@tencent.com
cuso4cdw@gmail.com, licm@smbu.edu.cn

Abstract

LLM agents have emerged as powerful systems for tackling multi-turn tasks by interleaving internal reasoning and external tool interactions. Agentic Reinforcement Learning has recently drawn significant research attention as a critical post-training paradigm to further refine these capabilities. In this paper, we present AT²PO (Agentic Turn-based Policy Optimization via Tree Search), a unified framework for multi-turn agentic RL that addresses three core challenges: limited exploration diversity, sparse credit assignment, and misaligned policy optimization. AT²PO introduces a turn-level tree structure that jointly enables *Entropy-Guided Tree Expansion* for strategic exploration and *Turn-wise Credit Assignment* for fine-grained reward propagation from sparse outcomes. Complementing this, we propose *Agentic Turn-based Policy Optimization*, a turn-level learning objective that aligns policy updates with the natural decision granularity of agentic interactions. ATPO is orthogonal to tree search and can be readily integrated into any multi-turn RL pipeline. Experiments across seven benchmarks demonstrate consistent improvements over the state-of-the-art baseline by up to 1.84 percentage points in average, with ablation studies validating the effectiveness of each component.

1 Introduction

Recent advances in large language models (LLMs) have catalyzed the development of autonomous agents capable of executing complex, multi-turn tasks through structured reasoning and tool use (Yao et al., 2023b; Schick et al., 2023; Shen et al., 2024). Reinforcement learning (RL) offers a natural framework for refining these agents by optimizing behavior toward task-level success. Inspired by recent progress in RL with Verifiable Rewards

* Equal contribution.

† Corresponding author.

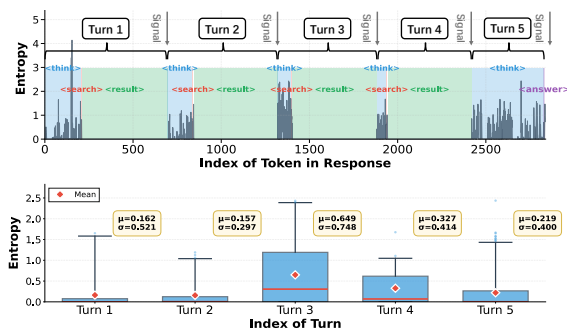


Figure 1: *Above*: Token-level entropy distribution in a search agentic rollout sample. A complete rollout sequence is sampled from multiple turns. *Below*: Turn-wise entropy box plot analysis. Substantial discrepancies exist among turn-level token groups.

(RLVR) which leverages programmatic feedback to guide training (Guo et al., 2025; Wang et al., 2025c; Wen et al., 2025; Lambert et al., 2025), Agentic RL extends this paradigm to interactive, multi-turn scenarios. In this context, agents act through turns of interleaved thinking and tool use, with tool responses forming the subsequent observations, creating a structured, turn-based interaction that requires specialized algorithmic designs. This emerging paradigm of agentic RL has drawn significant attention in the research community (Jin et al., 2025; Wu et al., 2025a; Dong et al., 2025a; Feng et al., 2025; Yan et al., 2025).

Despite its promise, Agentic RL still faces three fundamental challenges that hinder effective policy learning. First, agents struggle to generate high-quality, diverse trajectories under limited rollout budgets. Existing frameworks typically rely on chain-based generation or tree-based rollouts that expand nodes using random or heuristic criteria. These methods fail to strategically prioritize the exploration of high-uncertainty or high-potential turns, thereby limiting trajectory diversity and quality. Second, learning signals in multi-turn trajectories are inherently sparse. Rewards are typically available only at the end of a complete trajectory,

making it difficult to attribute credit to specific intermediate action steps. **Most critically, there remains a fundamental misalignment between the turn-based structure of agentic tasks and the flat optimization objectives used in current policy learning methods.** Current policy optimization algorithms inherited from RLVR treat the agent’s output as a flat token sequence by simply masking out the tool response tokens, and apply either token-level (Shao et al., 2024) or sequence-level (Zheng et al., 2025) importance sampling ratios. However, both schemes overlook a key structural property: the trajectory is generated through interleaved turns of agent reasoning and tool invocation, with new tool-generated tokens dynamically inserted between segments, as illustrated in Figure 1. Consequently, neither token-level nor sequence-level importance weighting properly accounts for the turn-wise sampling process. Moreover, even when step-wise learning signals are available, the unit of optimization mismatched with the natural unit of turn-level supervision. This structural mismatch prevents existing methods from capturing the step-wise logic of multi-turn decision-making, leading to unstable gradients and inefficient learning.

To overcome these challenges, we propose AT²PO (Agentic Turn-based Policy Optimization via Tree Search), a unified framework that integrates three synergistic components tailored to multi-turn agentic RL. 1) In the rollout phase, we introduce *Entropy-Guided Tree Expansion*, which adaptively grows the search tree from the most uncertain turns to maximize exploration efficiency. 2) In the rewarding phase, we develop a *Turn-wise Credit Assignment* mechanism that propagates sparse outcome rewards backward through the tree to compute fine-grained, per-turn value and advantage estimates. It enables precise credit allocation without auxiliary supervision. 3) In the training phase, we design *Agentic Turn-based Policy Optimization (ATPO)*, a novel policy learning algorithm that operates its importance sampling and clipping at the turn level, aligning the optimization objective with the structured turn-based decision process and significantly improving training stability and performance. Together, these components enable AT²PO to generate more diverse and high-quality rollouts, leverage sparse rewards more effectively, and optimize policies in a manner faithful to the multi-turn agentic paradigm. In conclusion, our contributions are summarized as follows:

- We propose a unified tree-based approach that jointly leverages Entropy-Guided Tree Expansion for strategic exploration and Turn-wise Credit Assignment for fine-grained reward propagation. The turn-level tree maximizes exploration diversity and supervision quality from sparse outcome rewards.
- We propose Agentic Turn-based Policy Optimization (ATPO), a turn-level policy update mechanism that is orthogonal to the tree-search framework and can be readily integrated into any multi-turn agentic RL pipeline. By operating at the natural granularity of multi-turn interactions, ATPO improves optimization stability and alignment compared to token-level or sequence-level baselines.
- Through extensive experiments on search agent benchmarks, we demonstrate that AT²PO consistently outperforms prior strong baselines, yielding gains of up to 1.84 percentage points in average over the state-of-the-art baseline across seven benchmarks. Ablation studies further validate the effectiveness of each component above.

2 Related Work

2.1 Reinforcement Learning for LLMs

RL has become a key approach for aligning LLMs with complex objectives and improving reasoning capabilities (Yu et al., 2025b; Zheng et al., 2025; Jin et al., 2025; Wang et al., 2025d). Traditional RLHF pipelines rely on Proximal Policy Optimization (PPO) (Schulman et al., 2017) with a learned reward model and value critic (Christiano et al., 2017; Stiennon et al., 2020). More recently, methods like RLVR leverage verifiable, programmatic rewards, demonstrating strong performance in settings with deterministic feedback (Guo et al., 2025). To simplify the architecture and reduce variance, a line of work replaces the critic with group-based baselines. GRPO (Guo et al., 2025) computes relative advantages within response groups, achieving strong results on long-form reasoning. DAPO (Yu et al., 2025b) improves scalability through decoupled clipping and dynamic sampling, while GSPO (Zheng et al., 2025) shifts clipping from the token to the sequence level for enhanced stability. Concurrent efforts further stabilize training via adaptive clipping strategies in the policy objec-

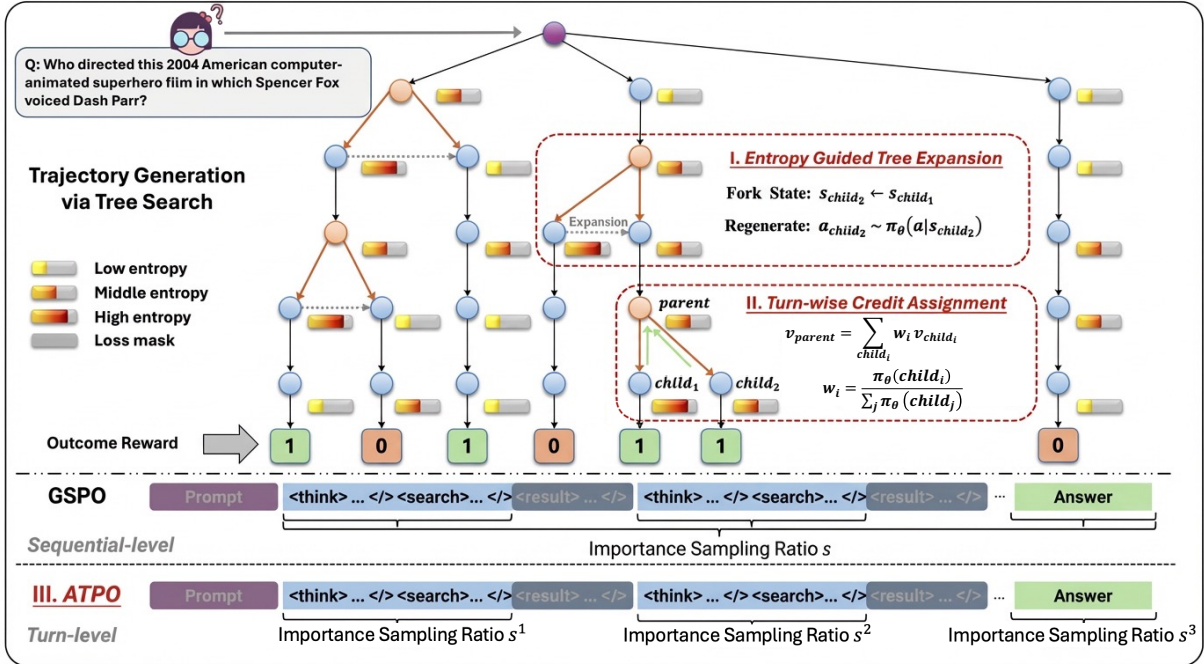


Figure 2: Overview of the **AT²PO** framework, including entropy-guided tree-structured rollout, turn-wise credit assignment for fine-grained supervision and turn-based policy optimization during reinforcement training.

tive (Chen et al., 2025; Wang et al., 2025b; Su et al., 2025b; Gao et al., 2025).

However, existing methods treat rollouts as flat sequences and do not explicitly account for the multi-turn nature of agentic interactions. In contrast, ATPO formulates the learning objective at the turn level, aligning policy updates with turn-wise feedback to improve training stability and quality.

2.2 Agentic Reinforcement Learning

Recent studies on Agentic RL focus primarily on trajectory data (Wu et al., 2025a; Li et al., 2025), algorithmic optimization (Dong et al., 2025b,a; Feng et al., 2025), memory management (Yan et al., 2025; Yu et al., 2025a), as well as multimodality (Geng et al., 2025; Wu et al., 2025b) and efficiency (Wang et al., 2025a). Algorithmically, several works focus on structured exploration. Search-R1 (Jin et al., 2025) integrates search tools into RL training, establishing a paradigm for knowledge-augmented agents. ARPO (Dong et al., 2025b) branches rollouts at high-entropy tokens to encourage exploration, while AEPO (Dong et al., 2025a) limits over-branching via entropy-aware clipping. GiGPO (Feng et al., 2025) introduces multi-granularity advantage estimation for finer credit assignment. Most closely related, Tree-GRPO (Ji et al., 2025) organizes rollouts into trees and computes group-relative advantages both within and across branches, extending GRPO to

agentic settings.

Even though methods like Tree-GRPO proposed to adopt a tree structure for more efficient rollout and trajectory collection, they do not explicitly prioritize expansion. In contrast, our entropy-guided tree expansion actively selects high-entropy nodes for branching, maximizing exploration efficiency. Coupled with ATPO our approach ensures coherent, diverse, and high-quality agentic behavior.

3 Preliminaries

An agentic framework interacts with external tools over multiple steps, which we model as a Markov Decision Process (MDP). At each time step t , the agent observes a state s_t and selects an action a_t , resulting in a transition to the next state s_{t+1} .

Following the widely adopted ReAct paradigm (Yao et al., 2023b), each action a_t comprises two components, including an internal chain-of-thought reasoning step and an executable operation. This operation is either a request to invoke an external tool or a final answer returned to terminate the interaction. Given a question in the input prompt, the agent iteratively samples actions $a_t \sim \pi_{\theta}(\cdot|s_t)$ according to its current policy π_{θ} , continuing until it emits the final answer at step $T - 1$. The resulting interaction trajectory is thus represented as $\tau = \{(s_t, a_t)\}_{t=0}^{T-1}$.

We assign a trajectory-level reward $R(\tau)$ based on the outcome, and aim to maximize the expected

return as follows:

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)], \quad (1)$$

Optimization is typically performed via policy gradient methods. In contrast to actor-critic approaches such as PPO (Schulman et al., 2017), our method follows group-based optimization (Guo et al., 2025), which leverages group-wise candidate sampling to implicitly estimate advantages without requiring an additional critic network.

4 Methods

In this section, we detail the methodology of AT²PO, which introduces specific enhancements across the rollout, rewarding, and training phases. First, during the rollout phase, AT²PO adopts a tree-based structure incorporating Entropy Guided Tree Expansion to enhance exploration diversity within a fixed inference budget. Second, in the rewarding phase, we leverage this tree topology to implement turn-wise credit assignment, enabling fine-grained supervision derived solely from outcome rewards. Finally, in the training phase, we design the turn-based policy optimization tailored for multi-turn agentic RL to enhance training stability and quality.

4.1 Entropy Guided Tree Expansion

During the rollout phase, tree-search strategies have demonstrated strong effectiveness in online agentic RL (Ji et al., 2025). In contrast to conventional RLVR settings where tree nodes correspond to specially selected individual tokens (Hou et al., 2025), multi-turn agentic tasks naturally lend themselves to a turn-level formulation. Each observation–action pair constitutes a single tree node. This clear contextual segmentation facilitates diverse action sampling (including reasoning steps and tool-use requests) from well-structured intermediate states.

Rather than relying on random or heuristic expansion strategies as in prior work (e.g., TreeGRPO (Ji et al., 2025)), we propose Entropy-Guided Tree Expansion, which iteratively expands the search tree from the most uncertain turns to promote diverse exploration. The full tree construction procedure is summarized in Algorithm 1.

Phase 1: Initialization. For each prompt $x_i \in \mathcal{P}$ in the training batch, we first generate M independent chain-based trajectories using the current policy π_θ . Each prompt x_i is associated with its own tree T_i initialized by adding these M trajectories as branches. All branches are rooted at a shared

Algorithm 1 Entropy Guided Tree Expansion

Require: Prompts \mathcal{P} , initial branches M , beam size K , expansion iterations L , policy π_θ , branching penalty coeff α

Ensure: Tree structures $\mathcal{T} = \{T_1, T_2, \dots, T_{|\mathcal{P}|}\}$

// Phase 1: Initialize Trees

```

1: for each prompt  $x_i \in \mathcal{P}$  do
2:   Initialize root node  $n_i$  with state  $s_{i,0} = x_i$ 
3:   for  $m = 1$  to  $M$  do
4:      $\tau_m \leftarrow \text{RolloutFrom}(s_{i,0}; \pi_\theta)$ 
5:     Add all nodes  $(s_{i,j}, a_{i,j})$  from  $\tau_m$  to  $T_i$ 
6:   end for
7: end for

```

// Phase 2: Node Sampling & Tree Expansion

```

8: for  $l = 1$  to  $L$  do
9:   for each tree  $T_i$  do
10:     $\mathcal{C} \leftarrow \text{get\_non\_leaf\_nodes}(T_i)$ 
11:     $s(n) = H_{\pi_\theta}(n) - \alpha |\text{children}(n.\text{parent})|$ 
12:     $\mathcal{S} \leftarrow \text{TopK}(\mathcal{C}, K, s)$ 
13:    for each selected node  $(s_{i,j}, a_{i,j}) \in \mathcal{S}$  do
14:       $\tau' \leftarrow \text{RolloutFrom}(s_{i,j}; \pi_\theta)$ 
15:      Add all nodes  $(s'_{i,j}, a'_{i,j})$  from  $\tau'$  to tree
16:       $T_i$  as new branch
17:    end for
18:  end for
19: return  $\mathcal{T}$ 

```

root node n_i storing the initial state $s_{i,0} = x_i$. All non-root nodes contain both the intermediate state and the corresponding action, i.e., the incremental rollout segment generated at that turn.

Phase 2.1: Sampling. We compute the policy entropy at each node to quantify decision uncertainty. Higher entropy indicates greater ambiguity in action selection, highlighting nodes that would benefit most from further exploration. Specifically, we use Monte-Carlo estimation to estimate the normalized entropy $H_{\pi_\theta}(n)$ of each node n :

$$\begin{aligned} H_{\pi_\theta}(n) &= \mathbb{E}_{y \sim \pi_\theta(\cdot|x)} [-\log \pi_\theta(y^k | x, y^{<k})] \\ &\approx \frac{1}{|y^k|} \sum_{y_t \in y^k} -\log \pi_\theta(y_t | x, y^{<k}) \end{aligned} \quad (2)$$

where y^k denotes the rollout segment generated in the k -th turn in a trajectory and stored at node n . We use the superscript y^k to indicate the subsequence corresponding to a single turn, and the subscript y_t to refer to the individual tokens within that subsequence.

Based on these entropy estimates, we select the

top K highest-entropy nodes in each tree and fork their states $s_{i,j}$ to initiate new branches. To encourage balanced exploration across the tree and prevent over-expansion from a single node, we introduce a branching penalty coefficient α that downweights repeatedly selected nodes.

Phase 2.2: Expansion. From each forked state $s_{i,j}$, we regenerate the remainder of the trajectory using the full context from the root to the fork point. The resulting continuation is appended as a new branch to the original tree.

We repeat Phase 2.1 and 2.2 for L iterations yielding a total of $M + LK$ leaf nodes per prompt. As the expected depth of the sampled nodes for expansion is half of the maximum length, given that the average token budget of an individual rollout as B , the total token budget for one search tree is $(M + LK/2)B$. Thus the tree-based approach generates a larger number of diverse rollouts under the same computational budget, enhancing both exploration efficiency and training signal quality.

4.2 Turn-wise Credit Assignment

During the rewarding stage, most existing RLVR algorithms solely relies on the final outcome reward (Shao et al., 2024) as ground-truth verification is typically only feasible upon completion of the full response. This sparse credit assignment obscures the contribution of individual actions to final outcomes, yielding weak or delayed learning signals that impede effective credit propagation in multi-turn agentic tasks.

Leveraging the search tree constructed during rollout, AT²PO enables fine-grained, turn-wise credit assignment. Specifically, we estimate the value V_n of each node n via Monte Carlo bootstrapping over its descendants,

$$V_n = \begin{cases} \hat{r}_n, & \text{if } n \text{ is leaf} \\ \sum_{c \in C(n)} w_c V_c, & \text{otherwise,} \end{cases} \quad (3)$$

$$w_c = \frac{\pi_\theta(c)}{\sum_{c' \in C(n)} \pi_\theta(c')},$$

where $C(n)$ denotes the set of child nodes of n , and $\hat{r}_n = \frac{r_n - \text{mean}\{r_n\}}{\text{std}\{r_n\}}$ is the normalized final outcome reward assigned via standard GRPO normalization only at leaves. This recursive aggregation effectively propagates the outcome-based rewards backward through the tree to each intermediate node.

Based on the node values, we compute turn-wise advantages $A_n = f(V_n)$ to serve as credit signals for policy updates, and assign them to all tokens

within. Several $f(\cdot)$ strategies are viable, including treating V_n directly as advantages, or other formulations proposed (Hou et al., 2025). We evaluate the effectiveness of different credit assignment schemes in Section 5.3.2.

4.3 Agentic Turn-based Policy Optimization

During training, we observe that existing agentic RL frameworks typically adopt optimization algorithms originally designed for conventional RLVR settings, such as GRPO in Tree-GRPO and GPPO (Su et al., 2025a) in AEPO. These methods treat the entire interaction as a monolithic sequence and apply standard policy gradients, without accounting for the unique structure of multi-turn agentic interactions.

Motivated by the success of recent sequence-level policy optimization methods like GSPO, which adjusts the learning objective into the sequence level to align with the final outcome reward, we propose Agentic Turn-Based Policy Optimization (ATPO) to better address the multi-turn nature of agentic tasks. Building on the fine-grained, turn-wise credit assignments derived from previous designs, we reformulate the policy optimization objective at the turn level for better alignment. Specifically, we introduce turn-wise importance sampling ratios and clipping mechanisms that operate per decision turn, enabling more precise and stable policy updates. This turn-aware formulation ensures that each reasoning or action step is (a) assigned importance ratios in a turn-consecutive manner and (b) optimized using its own turn-level learning signal.

$$\mathcal{J}_{ATPO}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} M_{i,t} \min \left(s_{i,t}^{turn}(\theta) \hat{A}_{i,t}, \text{clip} \left(s_{i,t}^{turn}(\theta), 1 - \varepsilon_l, 1 + \varepsilon_r \right) \hat{A}_{i,t} \right) \right], \quad (4)$$

$$s_{i,t}^{turn}(\theta) = \frac{\pi_\theta(y_{i,t}|x, y_{i,<t})}{\text{sg}[\pi_\theta(y_{i,t}|x, y_{i,<t})]} \cdot \text{sg} \left[\left(\frac{\pi_\theta(y_i^{k(t)}|x, y_i^{<k(t)})}{\pi_{\theta_{old}}(y_i^{k(t)}|x, y_i^{<k(t)})} \right)^{\frac{1}{|y_i^{k(t)}|}} \right], \quad (5)$$

where the individual token $y_{i,t}$ belongs to the turn $y_i^{k(t)}$, and $\text{sg}(\cdot)$ indicates the stop gradient operation. $M_{i,t}$ is a masking scheme that includes only

tokens generated by the inference engine, excluding those originating from tool responses.

The effectiveness of ATPO stems from its balanced design that bridges the gap between token-level and sequence-level policy optimization. First, in contrast to token-level methods such as GRPO, which suffer from high-variance of individual importance sampling ratios, ATPO operates at the turn level, significantly reducing clipped tokens and improving training stability. Second, unlike GSPO which facilitates its optimization objective and reward both on the coarse sequence-level, ATPO ensures that the optimization objective remains locally aligned with each consecutive turn-wise sampling step and with the fine-grained supervision provided by turn-based credit assignment.

Note that ATPO is orthogonal to the tree search framework, and serves as a plug-in design for any multi-turn agentic RL tasks. To better assess the effectiveness of turn-level optimization in ATPO, we also introduce a diagnostic metric called turn entropy \mathcal{H}_{turn} , which quantifies the degree of variation in policy updates across different turns within a trajectory. The details are presented in Appendix A.

5 Evaluation

5.1 Experiment Settings

Datasets. We evaluate our approach in a tool-integrated knowledge search setting. Following the instructions in Search-R1 (Jin et al., 2025), we implement a lightweight search engine and assess performance on seven widely adopted question answering benchmarks. These datasets are grouped into two categories: **Multi-Hop QA**—HotpotQA (Yang et al., 2018), 2WikiMultihopQA (Ho et al., 2020), MuSiQue (Trivedi et al., 2022), and Bamboogle (Press et al., 2023); and **Single-Hop QA**—Natural Questions (NQ) (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017), and PopQA (Mallen et al., 2022). All benchmarks are evaluated using Qwen3-4B, Qwen3-8B and Qwen2.5-7B models as the backbone, respectively. We use **Exact Match (EM)** as the primary evaluation metric across all datasets. Critically, our selection avoids reliance on proprietary or costly APIs and imposes minimal hardware requirements, ensuring that our evaluation is both reproducible and focused squarely on algorithmic advances.

Baselines. We compare AT²PO against the following widely-accepted RLVR baselines: GRPO

(Guo et al., 2025), DAPO (Yu et al., 2025b), GSPO (Zheng et al., 2025). We also compare AT²PO against two most recent agentic RL baselines: AEPO (Dong et al., 2025a) and Tree-GRPO (Ji et al., 2025). Note that we notice the performance of Tree-GRPO is highly unstable, even collapses at an early stage on Qwen-3 models. We provide further analysis on the possible cause in the Appendix C. We also report the base results by directly inferring via ReAct (Yao et al., 2023b).

Implementation Details. All experiments are conducted via respective RL recipes without any additional SFT phase. We use a training batch size of 64, a mini-batch size of 8 and the maximum response length of 6192. During rollout, we configure AT²PO with $M = 10$, $L = 2$ and $K = 6$, while other baselines except DAPO use a global rollout size of 16. These settings ensure a comparable total token budget across methods. The maximum tool usage is set to 6. The clipping thresholds for the ATPO objective are set to $3e-3$ and $4e-3$. All experiments are trained for 240 steps, and we compute the average accuracy across all evaluation samples. We select and report results from the checkpoint with the highest average score. Additional implementation specifics are provided in Section B.

5.2 Main Results of AT²PO

We report the performance of all baseline methods and AT²PO across different backbone models and benchmarks in Table 1.

Overall, AT²PO achieves the best results among most compared methods, yielding improvements of up to 1.84 percentage points in average over the state-of-the-art baseline across multi-hop and single-hop benchmarks. These results demonstrate the effectiveness and robustness of AT²PO in agentic RL. We note that AT²PO is slightly outperformed on certain single-hop QA datasets. A key reason is that DAPO employs a dynamic sampling strategy that utilizes up to three times the rollout budget compared to all other methods, including AT²PO. When evaluated under an equal computational budget to other baselines, AT²PO achieves the highest average performance across all settings.

Notably, AT²PO exhibits larger gains on multi-hop benchmarks than on single-hop ones. This stems from the fact that multi-hop tasks require significantly more interaction turns to retrieve the correct answer, allowing the turn-level design of AT²PO to fully demonstrate its advantage. As shown in the turn count distribution per sample

Method	Multi-Hop QA					Single-Hop QA			
	Hotpot	2wiki	Musiq	Bamb	Avg.	NQ	TriviaQA	PopQA	Avg.
Backbone Model: Qwen3-4B									
ReAct	30.42	32.92	12.83	44.80	30.01	26.75	53.53	35.34	41.31
+ GRPO	44.76	51.40	21.60	50.40	46.02	45.98	65.17	49.18	54.97
+ DAPO	45.95	51.81	21.68	51.20	46.65	47.50	65.84	51.03	56.33
+ GSPO	47.07	49.25	22.68	50.40	45.69	46.01	64.24	48.50	54.28
+ AEPO	46.36	51.78	23.47	50.40	46.95	45.71	64.66	50.13	55.20
+ AT ² PO (Ours)	49.44	52.99	24.80	56.80	48.81	47.90	65.32	51.81	56.44
Backbone Model: Qwen3-8B									
ReAct	20.66	19.05	9.56	37.60	18.66	21.16	41.81	27.37	32.19
+ GRPO	47.01	53.69	21.35	54.40	48.03	45.70	67.42	50.17	56.29
+ DAPO	49.64	53.91	24.05	56.00	49.40	51.99	69.02	51.90	58.53
+ GSPO	49.59	52.55	24.35	54.4	48.56	45.56	67.75	49.66	56.15
+ AEPO	49.17	52.97	24.01	54.40	48.62	49.92	68.31	51.77	57.94
+ AT ² PO (Ours)	51.37	53.97	26.51	56.00	50.15	51.33	69.51	52.26	58.82
Backbone Model: Qwen2.5-7B									
ReAct	2.85	1.94	0.58	4.00	2.10	4.34	10.67	9.32	9.23
+ GRPO	47.94	46.89	21.27	47.20	44.48	45.56	64.86	49.92	55.20
+ DAPO	47.50	47.93	21.27	44.00	44.91	52.24	65.00	50.01	56.08
+ GSPO	47.35	47.30	20.32	44.00	44.40	49.64	62.87	49.75	54.81
+ Tree-GRPO	42.39	42.01	20.15	42.40	39.79	47.56	62.69	44.75	52.04
+ AEPO	47.05	47.53	21.03	44.00	44.51	49.00	64.13	50.21	55.45
+ AT ² PO (Ours)	49.58	48.04	22.56	51.20	45.83	52.91	64.90	50.44	56.34

Table 1: Experiment results on three backbone models across seven datasets. The **bolded** values indicate the best result in comparisons. Our AT²PO outperforms existing methods in the majority of cases.

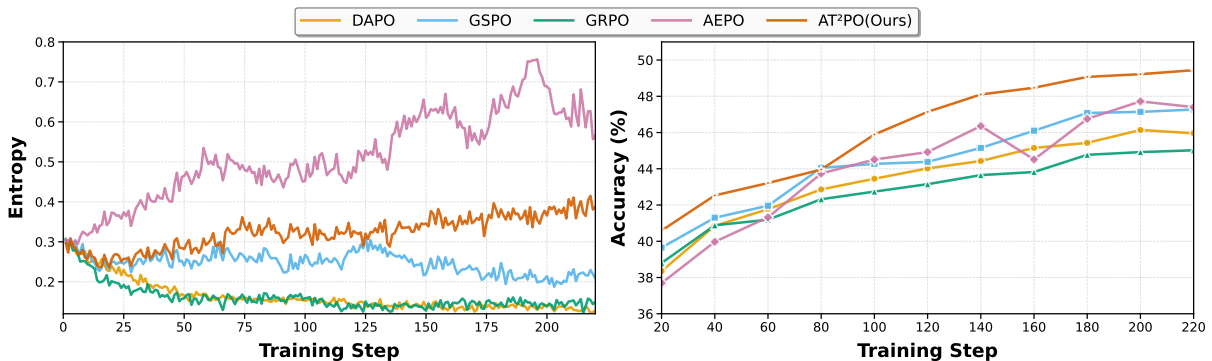


Figure 3: *Left*: The entropy comparison of different methods during training on multihop benchmark. *Right*: Performance comparison between all methods on HotpotQA dataset. Both are based on Qwen3-4B.

in the validation set in Figure 4, while 43.5% of samples in the single-hop benchmark require only a single tool call followed immediately by the answering turn, nearly all samples in the multi-hop benchmark necessitate at least two tool interactions. Consequently, AT²PO achieves state-of-the-art performance across all multi-hop settings, while still delivering strong results on single-hop benchmarks.

Furthermore, we analyze the training dynamics of AT²PO in detail. Figure 3 shows entropy evolution and validation accuracy across training steps. We observe that standard token-level optimization methods such as GRPO suffer from early entropy collapse, which severely limits their abil-

ity to explore high-quality trajectories. AEPO, which adapts a modified objective from GPPO (Su et al., 2025a) using a soft clipping strategy to preserve gradients for all tokens, avoids early collapse but exhibits gradual entropy divergence over long-horizon training. In contrast, AT²PO maintains the most stable entropy trajectory, effectively balancing exploration with gradient optimization stability.

5.3 Analysis

5.3.1 Ablation Study on AT²PO

Table 2 and Figure 8 present an ablation study on Qwen3-4B to assess the contribution of each component in AT²PO, following a progressive modular

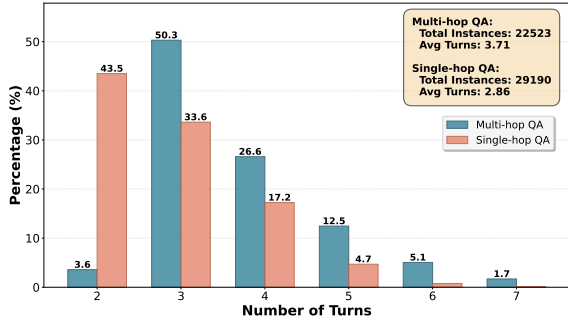


Figure 4: Turn count distribution per sample in the validation set across Multi-hop and Single-hop benchmarks with Qwen3-4B fine-tuned using AT²PO.

Model	Avg.
Qwen3-4B with ReAct	30.01
Random Tree Search with GRPO Loss	45.42
Random Tree Search with ATPO Loss	47.75
+ Entropy Guided Tree Expansion	48.33
+ Turn-wise Credit Assignment (AT ² PO)	48.81

Table 2: Ablation Study on progressive modular designs applied to AT²PO on multi-hop benchmarks.

where modules are incrementally added. This setup enables us to isolate the effect of each design choice under a fixed backbone and consistent task suite.

Starting from a base configuration of standard GRPO with random tree search, replacing it with ATPO yields substantial gains across all multi-hop benchmarks. This highlights that the turn-wise importance sampling and clipping provides a significantly finer-grained optimization objective than vanilla GRPO in multi-turn agentic RL. Next, integrating Entropy-Guided Tree Expansion leads to a further clear performance improvement. This indicates that entropy-aware branching facilitates more effective exploration during tree-based rollouts, enabling the agent to allocate its limited search budget toward more informative decision points. Finally, adding Turn-wise Credit Assignment delivers additional gains and achieves the best overall performance. This demonstrates that providing fine-grained, step-level learning signals to individual turns effectively mitigates the credit assignment problem inherent in sparse-reward agentic tasks.

5.3.2 Scheme Analysis of Turn-wise Credit Assignment

Table 3 and Table 5 compare different formulations of turn-wise credit assignment by varying both the advantage definition and the value aggregation strategy. We consider three strategies for estimating the node value V_n : (a) child-weighted

Advantage $A(n)$	Value $V(n)$	Avg.
V_n	Child Weighted	48.81
V_n	Child Mean	48.43
V_n	Leaf Mean	48.52
$L_n = V_n - V_{n_{\text{parent}}}$	Child Weighted	45.80
$L_n = V_n - V_{n_{\text{parent}}}$	Leaf Mean	45.61
$G_n = V_n - V_{n_{\text{root}}}$	Child Weighted	48.19
$G_n = V_n - V_{n_{\text{root}}}$	Leaf Mean	46.98
$L_n + G_n$	Leaf Mean	44.18

Table 3: Ablation on the effectiveness of different turn-wise credit assignment choices.

propagation as defined in in Eq (3), (b) mean of all child values and (c) mean of all leaf values following TreeRL (Hou et al., 2025). As for advantage computation $f(\cdot)$, our main experiments adopt the simplest formulation $A_n = V_n$. We further compare this with several difference-based alternatives: (a) the TD advantage $A_n = L_n = V_n - V_{n_{\text{parent}}}$, (b) the global advantage $A_n = G_n = V_n - V_{n_{\text{root}}}$ and (c) the combined signal $A_n = L_n + G_n$ as used in prior work (Hou et al., 2025).

Overall, directly setting $A_n = V_n$ yields best performance across all datasets. Under this formulation, child-weighted aggregation achieves the best results, while child-mean and leaf-mean aggregation remain competitive but slightly inferior. In contrast, difference-based formulations lead to noticeably lower performance. Surprisingly, the combined signal, despite its demonstrated effectiveness in conventional RLVR settings (Yao et al., 2023a), performs worst among all variants. In the agentic RL, directly using V_n as the advantage provides a cleaner and more compatible learning signal.

6 Conclusion

In this paper, we present AT²PO, a unified framework for multi-turn agentic RL. Through Entropy-Guided Tree Expansion, we enable strategic exploration at uncertain decision points; via Turn-wise Credit Assignment, we propagate sparse outcome rewards to derive fine-grained supervision signals; and with Agentic Turn-based Policy Optimization, we align policy updates with the natural turn-based structure. Extensive experiments across seven benchmarks demonstrate consistent improvements over strong baselines, validating the effectiveness of each component. Our framework offers a promising approach to training LLM agents for complex multi-turn reasoning and tool-use tasks.

7 Limitations

While AT²PO achieves strong empirical performance, its tree-based expansion introduces additional computational overhead compared to linear rollouts particularly when sufficient computational resources are available, as the expansion process requires multiple sequential iterations. Future work will explore more efficient rollout strategies through enhanced parallelization. Additionally, we aim to evaluate AT²PO in a broader range of agentic environments to further assess its generalization and robustness.

References

- Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang, Bo Fei, Bo Yang, Boji Shan, Changqing Yu, Chao Wang, Cheng Zhu, and 1 others. 2025. Minimax-m1: Scaling test-time compute efficiently with lightning attention. *arXiv preprint arXiv:2506.13585*.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Guanting Dong, Licheng Bao, Zhongyuan Wang, Kangzhi Zhao, Xiaoxi Li, Jiajie Jin, Jinghan Yang, Hangyu Mao, Fuzheng Zhang, Kun Gai, and 1 others. 2025a. Agentic entropy-balanced policy optimization. *arXiv preprint arXiv:2510.14545*.
- Guanting Dong, Hangyu Mao, Kai Ma, Licheng Bao, Yifei Chen, Zhongyuan Wang, Zhongxia Chen, Jiazhen Du, Huiyang Wang, Fuzheng Zhang, and 1 others. 2025b. Agentic reinforced policy optimization. *arXiv preprint arXiv:2507.19849*.
- Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. 2025. Group-in-group policy optimization for llm agent training. *arXiv preprint arXiv:2505.10978*.
- Chang Gao, Chujie Zheng, Xiong-Hui Chen, Kai Dang, Shixuan Liu, Bowen Yu, An Yang, Shuai Bai, Jingren Zhou, and Junyang Lin. 2025. [Soft adaptive policy optimization](#). *Preprint*, arXiv:2511.20347.
- Xinyu Geng, Peng Xia, Zhen Zhang, Xinyu Wang, Qiuchen Wang, Ruixue Ding, Chenxi Wang, Jialong Wu, Yida Zhao, Kuan Li, and 1 others. 2025. Webwatcher: Breaking new frontier of vision-language deep research agent. *arXiv preprint arXiv:2508.05748*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. [Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Zhenyu Hou, Ziniu Hu, Yujiang Li, Rui Lu, Jie Tang, and Yuxiao Dong. 2025. [Treeerl: Llm reinforcement learning with on-policy tree search](#). *Preprint*, arXiv:2506.11902.
- Yuxiang Ji, Ziyu Ma, Yong Wang, Guanhua Chen, Xi-angxiang Chu, and Liaoni Wu. 2025. [Tree search for llm agent reinforcement learning](#). *Preprint*, arXiv:2509.21240.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, and 4 others. 2025. [Tulu 3: Pushing frontiers in open language model post-training](#). *Preprint*, arXiv:2411.15124.
- Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, and 1 others. 2025. Websailor: Navigating super-human reasoning for web agent. *arXiv preprint arXiv:2507.02592*.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Hannaneh Hajishirzi, and Daniel Khashabi. 2022. When not to trust language models: Investigating effectiveness and limitations of parametric and non-parametric memories. *arXiv preprint*.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. [Measuring](#)

- and narrowing the compositionality gap in language models. *Preprint*, arXiv:2210.03350.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. [Qwen2.5 technical report](#). *Preprint*, arXiv:2412.15115.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. [Toolformer: Language models can teach themselves to use tools](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. [Proximal policy optimization algorithms](#). *Preprint*, arXiv:1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *Preprint*, arXiv:2402.03300.
- Yuanhao Shen, Xiaodan Zhu, and Lei Chen. 2024. Smartcal: An approach to self-aware tool-use evaluation and calibration. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 774–789.
- Guangming Sheng, Chi Zhang, Zilinfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:2409.19256*.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008–3021.
- Zhenpeng Su, Leiyu Pan, Xue Bai, Dening Liu, Guanting Dong, Jiaming Huang, Wenping Hu, Fuzheng Zhang, Kun Gai, and Guorui Zhou. 2025a. Klear-reasoner: Advancing reasoning capability via gradient-preserving clipping policy optimization. *arXiv preprint arXiv:2508.07629*.
- Zhenpeng Su, Leiyu Pan, Minxuan Lv, Yuntao Li, Wenping Hu, Fuzheng Zhang, Kun Gai, and Guorui Zhou. 2025b. Ce-gppo: Coordinating entropy via gradient-preserving clipping policy optimization in reinforcement learning. *arXiv preprint arXiv:2509.20712*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. [MuSiQue: Multi-hop questions via single-hop question composition](#). *Transactions of the Association for Computational Linguistics*, 10:539–554.
- vLLM Team. 2025. No more retokenization drift: Returning token ids via the openai compatible api matters in agent rl. <https://blog.vllm.ai/2025/10/22/agent-lightning.html>. Accessed: 2026-01-05.
- Hongru Wang, Cheng Qian, Wanjun Zhong, Xiusi Chen, Jiahao Qiu, Shijue Huang, Bowen Jin, Mengdi Wang, Kam-Fai Wong, and Heng Ji. 2025a. Otc: Optimal tool calls via reinforcement learning. *arXiv e-prints*, pages arXiv–2504.
- Jiakang Wang, Runze Liu, Lei Lin, Wenping Hu, Xiu Li, Fuzheng Zhang, Guorui Zhou, and Kun Gai. 2025b. Aspo: Asymmetric importance sampling policy optimization. *arXiv preprint arXiv:2510.06062*.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*.
- Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Liyuan Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, Weizhu Chen, Shuhang Wang, Simon Shaolei Du, and Yelong Shen. 2025c. [Reinforcement learning for reasoning in large language models with one training example](#). *Preprint*, arXiv:2504.20571.
- Zihan Wang, Kangrui Wang, Qineng Wang, Pingyue Zhang, Linjie Li, Zhengyuan Yang, Xing Jin, Kefan Yu, Minh Nhat Nguyen, Licheng Liu, and 1 others. 2025d. Ragen: Understanding self-evolution in llm agents via multi-turn reinforcement learning. *arXiv preprint arXiv:2504.20073*.
- Xumeng Wen, Zihan Liu, Shun Zheng, Shengyu Ye, Zhirong Wu, Yang Wang, Zhijian Xu, Xiao Liang, Junjie Li, Ziming Miao, Jiang Bian, and Mao Yang. 2025. [Reinforcement learning with verifiable rewards implicitly incentivizes correct reasoning in base llms](#). *Preprint*, arXiv:2506.14245.
- Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Gang Fu, Yong Jiang, and 1 others. 2025a. Webdancer: Towards autonomous information seeking agency. *arXiv preprint arXiv:2505.22648*.
- Jinming Wu, Zihao Deng, Wei Li, Yiding Liu, Bo You, Bo Li, Zejun Ma, and Ziwei Liu. 2025b. Mmsearch-r1: Incentivizing llms to search. *arXiv preprint arXiv:2506.20670*.
- Sikuan Yan, Xiufeng Yang, Zuchao Huang, Ercong Nie, Zifeng Ding, Zonggen Li, Xiaowen Ma, Kristian Kersting, Jeff Z Pan, Hinrich Schütze, and 1 others. 2025. Memory-r1: Enhancing large language model agents to manage and utilize memories via reinforcement learning. *arXiv preprint arXiv:2508.19828*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,

- Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023b. [React: Synergizing reasoning and acting in language models](#). In *The Eleventh International Conference on Learning Representations*.
- Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie Chen, Weinan Dai, Qiyong Yu, Ya-Qin Zhang, Wei-Ying Ma, Jingjing Liu, Mingxuan Wang, and 1 others. 2025a. Memagent: Reshaping long-context llm with multi-conv rl-based memory agent. *arXiv preprint arXiv:2507.02259*.
- Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, YuYue, Weinan Dai, Tiantian Fan, Gao-hong Liu, Juncai Liu, LingJun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, and 17 others. 2025b. [DAPO: An open-source LLM reinforcement learning system at scale](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Chi Zhang, Guangming Sheng, Siyao Liu, Jiahao Li, Ziyuan Feng, Zherui Liu, Xin Liu, Xiaoying Jia, Yanghua Peng, Haibin Lin, and 1 others. 2024. A framework for training large language models for code generation via proximal policy optimization. In *NL2Code Workshop of ACM KDD*.
- Chujie Zheng, Pei Ke, Zheng Zhang, and Minlie Huang. 2023. Click: Controllable text generation with sequence likelihood contrastive learning. *arXiv preprint arXiv:2306.03350*.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, Jingren Zhou, and Junyang Lin. 2025. [Group sequence policy optimization](#). *Preprint*, arXiv:2507.18071.

A Analysis of ATPO Objective

In this section, we clarify the distinctions among token-level, sequence-level, and turn-level policy optimization objectives, highlighting how each approach leverages different granularities of policy updates.

Token-Level Optimization. GRPO (Shao et al., 2024) computes the relative advantage of each response within a group and optimizes the following token-level objective:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} M_{i,t} \min \left(r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip} \left(r_{i,t}(\theta), 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{i,t} \right) \right], \quad (6)$$

where the importance sampling ratio is defined at the token level as

$$r_{i,t}(\theta) = \frac{\pi_{\theta}(y_{i,t}|x, y_{i,<t})}{\pi_{\theta_{old}}(y_{i,t}|x, y_{i,<t})}. \quad (7)$$

Notably, All tokens in a sample share the same advantage $\hat{A}_{i,t} = \hat{A}_i$.

Sequence-Level Optimization. In contrast, GSPO (Zheng et al., 2025) adopts a sequence-level perspective by defining the importance ratio based on the entire response likelihood:

$$\mathcal{J}_{GSPO}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} M_{i,t} \min \left(s_{i,t}(\theta) \hat{A}_{i,t}, \text{clip} \left(s_{i,t}(\theta), 1 - \varepsilon_l, 1 + \varepsilon_r \right) \hat{A}_{i,t} \right) \right], \quad (8)$$

with the sequence-aware importance ratio given by

$$s_{i,t}(\theta) = \text{sg} \left[\left(\frac{\pi_{\theta}(y_i|x)}{\pi_{\theta_{old}}(y_i|x)} \right)^{\frac{1}{|y_i|}} \cdot \frac{\pi_{\theta}(y_{i,t}|x, y_{i,<t})}{\text{sg}[\pi_{\theta}(y_{i,t}|x, y_{i,<t})]} \right]. \quad (9)$$

Here, $\text{sg}[\cdot]$ denotes the stop-gradient operator. By anchoring the importance weight to the full-sequence likelihood (following (Zheng et al.,

2023)), GSPO ensures that clipping is applied coherently across all tokens in a response, reflecting an all-or-nothing treatment of off-policy sequences.

Turn-Level Optimization. Our proposed Adaptive Turn-level Policy Optimization (ATPO) introduces a turn-aware objective that balances fine-grained control with sequence coherence. The objective is:

$$\mathcal{J}_{ATPO}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} M_{i,t} \min \left(s_{i,t}^{\text{turn}}(\theta) \hat{A}_{i,t}, \text{clip} \left(s_{i,t}^{\text{turn}}(\theta), 1 - \varepsilon_l, 1 + \varepsilon_r \right) \hat{A}_{i,t} \right) \right], \quad (10)$$

where the turn-level importance ratio is defined as

$$s_{i,t}^{\text{turn}}(\theta) = \frac{\pi_{\theta}(y_{i,t}|x, y_{i,<t})}{\text{sg}[\pi_{\theta}(y_{i,t}|x, y_{i,<t})]} \cdot \text{sg} \left[\left(\frac{\pi_{\theta}(y_i^{k(t)}|x, y_i^{<k(t)})}{\pi_{\theta_{old}}(y_i^{k(t)}|x, y_i^{<k(t)})} \right)^{\frac{1}{|y_i^{k(t)}|}} \right]. \quad (11)$$

Here $y_{i,t}$ belongs to the $k(t)$ -th response turn $y_i^{k(t)}$. This formulation enables selective gradient propagation: if tokens from a particular turn are highly off-policy, their contribution to the gradient can be suppressed without discarding updates from other more on-policy turns. The gradient is derived as follows,

$$\nabla_{\theta} \mathcal{J}_{ATPO}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} M_{i,t} s_{i,t}^{\text{turn}}(\theta) \hat{A}_{i,t} \cdot \nabla_{\theta} \log \pi_{\theta}(y_{i,t}|x, y_{i,<t}) \right] \quad (12)$$

To better assess the effectiveness of turn-level optimization in ATPO, we also introduce a diagnostic metric called turn entropy $\mathcal{H}_{\text{turn}}$, which quantifies the degree of variation in policy updates across different turns within a trajectory:

$$\mathcal{H}_{\text{turn}} = \frac{1}{B} \sum_{i=1}^B \frac{-\sum_{j=1}^{N_i} p_{i,j} \log(p_{i,j})}{\log(N_i)}, \quad (13)$$

$$p_{i,j} = \frac{\exp\left(\frac{KL_{i,j}}{KL_{i,seq}}\right)}{\sum_{k=1}^{N_i} \exp\left(\frac{KL_{i,k}}{KL_{i,seq}}\right)},$$

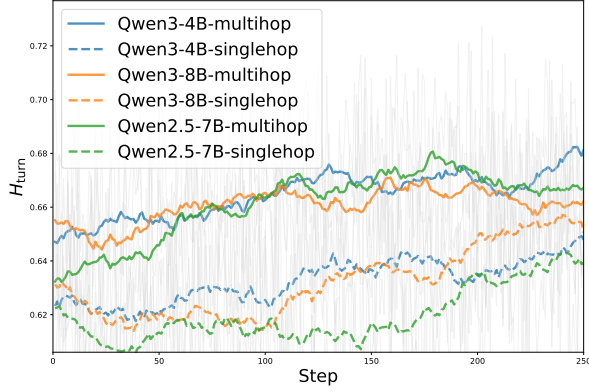


Figure 5: H_{turn} alongside training steps of AT²PO in all experiment settings.

where N_i is number of turns of the i -th sample in batch, $KL_{i,j}$ is the KL divergence between the old and new policies at turn j , and $KL_{i,seq}$ is the total KL divergence of the full sequence.

We normalize \mathcal{H}_{turn} to $[0, 1]$. A special value of $\mathcal{H}_{turn} = 1$ indicates uniform policy updates across all turns, while $\mathcal{H}_{turn} = 0$ occurs when every sample contains only one turn. Conversely, smaller values signify greater heterogeneity in turn-level updates, highlighting scenarios where turn-aware optimization is particularly beneficial.

Figure 5 plots H_{turn} throughout training across all experiment settings. The entropy stabilizes around 0.66 for multi-hop benchmarks and 0.62 for single-hop benchmarks, confirming significant variation in turn-level policy update and justifying the need for turn-granular optimization. Interestingly, we find that H_{turn} gradually increases during training across all settings, suggesting that as the policy converges, update magnitudes become more balanced across turns, which is a sign of stabilized and high-quality learning.

B Implementation Details

B.1 Reward Design

We adopt a binary rule-based reward function to obtain the final outcome reward in our agentic RL framework, focusing on both answer correctness and format completeness of LLM’s output. To be detailed, we follow the reward modeling setting of Search-R1 (Jin et al., 2025) which utilizes the EM-score as primary reward. Meanwhile, we also introduce the appropriate format constraint in this process.

Exact Match Reward Let \hat{y} denote the agent’s predicted final answer and y^* the ground-truth an-

swer. We leverage an **Exact Match (EM)** score as the primary metric of answer correctness:

$$r_{EM}(\hat{y}, y^*) = \begin{cases} 1, & \text{if } \hat{y} = y^* \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

This strict binary reward provides an unambiguous reward signal, encouraging the agentic LLM to produce fully correct answers rather than partially correct or ambiguous outputs. Compared to continuous value rewards like F1-score, it has a higher upper bound for reinforcement training.

Format Constraint In addition to the correctness of answer, we enforce a format validation constraint on the model output as well. Specifically, the generated response is necessary to contain both a reasoning trace and a final answer segment, explicitly delimited by the tags `<think>...</think>` and `<answer>...</answer>`, respectively. Besides, The final answer used for matching needs to be enclosed within `\boxed{\}` (within `<answer>...</answer>`). Formally, we define a format validation function as follows:

$$\mathbb{I}_{\text{format}} = \begin{cases} 1, & \text{if both tags are present} \\ 0, & \text{otherwise} \end{cases} \quad (15)$$

If the response violates the required format, the reward is set to zero regardless of the predicted answer correctness. This plays an important role for improving the model’s ability to follow user’s instructions.

Final Reward Definition Based on the above, the overall reward r is defined as the synthesis of exact match correctness and format constraint:

$$r = \begin{cases} r_{EM}(\hat{y}, y^*), & \text{if } \mathbb{I}_{\text{format}} = 1 \\ -1, & \text{otherwise} \end{cases} \quad (16)$$

Consequently, the agentic LLM derives a reward of 1 if and only if it outputs a correctly formatted output whose final answer exactly matches the ground truth.

B.2 Prompt Template

We present our prompt template as Figure 6. We follow the similar setting in (Jin et al., 2025) and (Dong et al., 2025a) with different tags containing corresponding content. Specifically, the model outputs its intermediate reasoning within the `<think></think>` tags and makes a search request by emitting the `<search></search>` tag, which corresponds to an action in the agent–environment interaction process. The content returned by the

Prompt Template.

You are a helpful assistant that can solve the given question step by step with the help of the wikipedia search tool. Given a question, you need to first think about the reasoning process in the mind and then provide the answer. During thinking, you can invoke the wikipedia search tool to search for fact information about specific topics if needed. You can search as many times as your want. The reasoning process and answer are enclosed within `<think>` `</think>` and `<answer>` `</answer>` tags respectively, and the search query and result are enclosed within `<search>` `</search>` and `<result>` `</result>` tags respectively. For example, `<think>` This is the reasoning process. `</think>` `<search>` search query here `</search>` `<result>` search result here `</result>` `<think>` This is the reasoning process. `</think>` `<answer>` The final answer is [`\boxed{answer here}`] `</answer>`. In the last part of the answer, the final exact answer is enclosed within `\boxed{}` with latex format. Question:

Figure 6: The prompt template.

search engine is wrapped in `<result>``</result>` tags and embedded into the response context. Finally, after completing the reasoning process, the model outputs the content of answer within the `<answer>``</answer>` tags. Noting that we exact the final answer from the `\boxed{}` in response to achieve more accurate EM score matching.

B.3 Datasets

We introduce two categories of benchmarks adopted in our experiments in this section.

Multi-Hop QA. To evaluate the multi-turn tool calling and complex reasoning capabilities, we adopt four datasets that require reasoning over multiple steps. **HotpotQA** (Yang et al., 2018) is a large-scale Wikipedia-based benchmark with strong supervision for supporting facts, widely used to evaluate explainable multi-hop reasoning. **2WikiMultiHopQA** (Ho et al., 2020) consists of Wikipedia text and Wikidata triples, which is a comprehensive benchmark with multi-hop dependency between entities. **Musique** (Trivedi et al., 2022) has 25k 2-4 hop questions synthesized via controlled composition of single-hop content. **Bamboogle** (Press et al., 2023) is a tiny dataset consisting of complex compositional QA questions, which benefits the evaluation of the effectiveness and stability of agentic RL methods.

Single-Hop QA. We further adopt three single-hop QA benchmarks. **Natural Questions (NQ)** (Kwiatkowski et al., 2019) is widely used to assess retrieval-augmented generation, which includes large quantity of user queries from Wikipedia. **TriviaQA** (Joshi et al., 2017) contains questions

with syntactic and lexical variability between questions and corresponding answer-evidence sentences. **PopQA** (Mallen et al., 2022) is an entity-centric open-domain QA dataset, designed to explore the interaction between agentic retrieval and parametric memorization.

B.4 Baseline Settings

In Table 4, we exhibit common hyperparameter settings used for baselines in main experiment. We further detail the baselines selected in our main experiments and illustrate the corresponding specially adjusted experiment settings.

- **ReAct** (Yao et al., 2023b): A prompting paradigm to enable procedural planning and dynamic function calls(tool calls), which is training-free.
- **GRPO** (Guo et al., 2025): A group-based on-policy RL optimization method, which generates n sets of candidates by multiple sampling for each prompt to derive the relative advantages, thereby replacing the traditional critic used in PPO. We set the coefficient to 0.001 for KL constraint. The clip ratio is set to 0.2.
- **DAPO** (Yu et al., 2025b): An extension of GRPO that incorporates a decoupled clipping mechanism and dynamic sampling to stabilize policy updates. Following the recommended configuration, we set `clip_ratio_low` and `clip_ratio_high` to 0.2 and 0.28, respectively, to enable a wider clipping range. Additionally, we configure the `overlong_buffer`

Config	Value
optimizer	AdamW
learning rate	1e-6
clip_ratio	0.2
total training steps	240
training batch size	64
PPO mini batch size	8
rollout_n	16
max prompt length	2000
max response length	6192
max tool-calls	6
reward metrics	EM
retriever	local wiki
top-K retrieval passages	3

Table 4: General hyperparameters of baseline methods in main experiments.

with a capacity of 2000 (penalty factor 1.0) and activate the dynamic sampling mechanism. Consistent with DAPO’s default setup, the generation batch size is set to three times the training batch size.

- **GSPO** (Zheng et al., 2025): A sequence-level variant of GRPO that defines importance ratios and clipping based on full response likelihoods, trading fine-grained token-level control for improved sequence-level stability. Analogous to DAPO, we adopt asymmetric clipping with $\text{clip_ratio_low} = 3 \times 10^{-4}$ and $\text{clip_ratio_high} = 4 \times 10^{-4}$ to support precise sequence-level clipping.
- **AEPO** (Dong et al., 2025a): A state-of-the-art agentic RL method that introduces entropy-balanced rollout scheduling and entropy-aware clipping to mitigate over-branching and gradient collapse in tree-based rollouts. In our experiments, we set initial_rollouts to 8. Additional hyperparameters include $\text{beam_size} = 2$, $\text{branch_probability} = 0.5$, and $\text{entropy_weight} = 0.2$.
- **Tree-GRPO** (Ji et al., 2025): A tree-based agentic RL framework that integrates GRPO with tree search to enable fine-grained credit assignment across branching trajectories. We replicate the original experiment configuration as reported in the paper.

B.5 Search Tool Settings

Follow the settings of search tool in Search-R1 (Jin et al., 2025), we utilize the wikipedia as the cor-

pus for search engine with e5-base-v2 (Wang et al., 2022) as the retriever model. The retrieval knowledge base contains 21 million data entries obtained from Wikipedia. After the agentic LLM calls the wiki search tool, the search engine will return the top-k most relevant knowledge data determined by the retrieval model at each turn.

B.6 Hardware and Artifacts

All experiments are conducted on a system equipped with 8 NVIDIA H20 GPUs. As the backbone models, we employ Qwen3-4B, Qwen3-8B, and Qwen2.5-7B (Yang et al., 2025; Qwen et al., 2025), chosen for their strong foundational reasoning capabilities in agentic RL training. Our implementation is built on the VeRL framework (Sheng et al., 2024; Zhang et al., 2024), a mature and stable RL infrastructure that can be easily integrated with agentic RL training features.

C Analysis on Training Stability

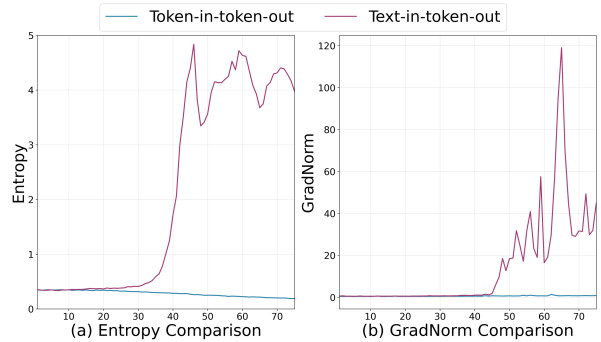


Figure 7: Training comparison between the scheme of token-in-token-out and text-in-token-out on Qwen3-4B with multi-hop training set.

During our replication of Tree-GRPO (Ji et al., 2025) using the open-sourced implementation, we observed significant training instability across multiple runs, often leading to complete collapse at very early stages when using Qwen3-based models. We trace this issue primarily to a phenomenon known as Retokenization Drift (vLLM Team, 2025).

In agentic RL rollouts in the Tree-GRPO implementation, intermediate responses are first detokenized into text strings to enable tool invocation and parsing. Tool outputs in string form are then concatenated to the dialogue history, and the entire sequence is retokenized before being used for policy updates. Despite semantic equivalence at the string level, this text-in-token-out pipeline can

Advantage $A(n)$	Value $V(n)$	Hotpot	2wiki	Musiq	Bamb	Avg.
$V(n)$	Child Weighted	49.44	52.99	24.80	54.4	48.81
$V(n)$	Child Mean	47.95	53.45	23.56	52.8	48.43
$V(n)$	Leaf Mean	49.10	52.76	24.18	54.4	48.52
$L(n) = V(n) - V(n_{\text{parent}})$	Child Weighted	44.73	51.38	19.73	51.2	45.80
$L(n) = V(n) - V(n_{\text{parent}})$	Leaf Mean	45.36	50.53	20.40	52.8	45.61
$G(n) = V(n) - V(n_{\text{root}})$	Child Weighted	48.45	52.87	22.73	54.4	48.19
$G(n) = V(n) - V(n_{\text{root}})$	Leaf Mean	47.67	50.96	23.68	52.8	46.98
$L(n) + G(n)$	Leaf Mean	45.35	48.24	19.07	52.0	44.18

Table 5: Ablation on the effectiveness of different turn-wise credit assignment schemes.

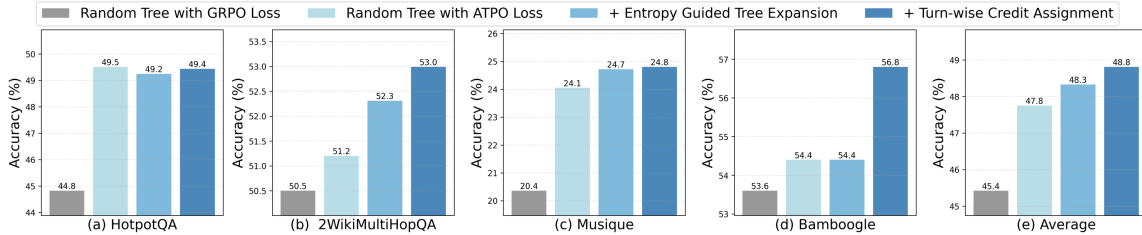


Figure 8: Ablation Study on progressive modular designs applied to AT²PO.

yield different token sequences before and after retokenization due to tokenizer nondeterminism, introducing critical inconsistencies between the sampled and training trajectories.

To address this, we adopt a token-in–token-out approach: we preserve the exact token IDs generated at each turn during rollout and use them directly in training, bypassing intermediate detokenization and retokenization. As illustrated in Figure 7, this modification eliminates retokenization-induced distributional shifts. Empirically, the revised pipeline stabilizes training—both policy entropy and gradient norm remain consistent throughout training, with no signs of collapse.

While this fix resolves a key source of instability, other challenges in stabilizing agentic RL training persist, including consistency between training and inference, etc. These remain important directions for future work beyond the scope of this paper.

D Detailed Ablations

D.1 Ablation Study on AT²PO

Table 2 and Figure 8 present an ablation study on Qwen3-4B to assess the contribution of each component in AT²PO, following a progressive modular where modules are incrementally added. This setup enables us to isolate the effect of each design choice under a fixed backbone and consistent task suite.

Starting from a base configuration that uses the standard GRPO loss with random tree search, replacing it with our proposed ATPO loss yields

substantial gains across all multi-hop benchmarks. This highlights that the turn-wise importance sampling and clipping mechanism provides a significantly finer-grained optimization objective than vanilla GRPO in multi-turn agentic RL. Next, integrating Entropy-Guided Tree Expansion leads to a further clear performance improvement. This indicates that entropy-aware branching facilitates more effective exploration during tree-based rollouts, enabling the agent to allocate its limited search budget toward more informative and uncertain decision points. Finally, adding Turn-wise Credit Assignment—completing the full AT²PO framework—delivers additional gains and achieves the best overall performance. This demonstrates that providing fine-grained, step-level learning signals to individual turns effectively mitigates the credit assignment problem inherent in sparse-reward, long-horizon agentic tasks.

D.2 Scheme Analysis of Turn-wise Credit Assignment

Table 3 and Table 5 compare different formulations of turn-wise credit assignment by varying both the advantage definition and the value aggregation strategy. We consider three strategies for estimating the node value V_n : (a) child-weighted propagation as defined in in Eq (3), (b) mean of all child values and (c) mean of all leaf values following TreeRL (Hou et al., 2025). As for advantage computation $f(\cdot)$, our main experiments adopt the

simplest formulation $A_n = V_n$. We further compare this with several difference-based alternatives: (a) the TD advantage $A_n = L_n = V_n - V_{n_{\text{parent}}}$, (b) the global advantage $A_n = G_n = V_n - V_{n_{\text{root}}}$ and (c) the combined signal $A_n = L_n + G_n$ as used in prior work (Hou et al., 2025).

Overall, directly setting the advantage equal to the node value yields consistently strong performance across all datasets. Under this formulation, child-weighted aggregation achieves the best results, while child-mean and leaf-mean aggregation remain competitive but slightly inferior. In contrast, difference-based formulations lead to noticeably lower performance. Surprisingly, the combined signal, despite its demonstrated effectiveness in conventional RLVR settings (Yao et al., 2023a), performs worst among all variants. In the agentic RL scenario, directly using V_n as the advantage provides a cleaner and more compatible learning signal.

E Case Study

We conduct the case study which is presented through Table 6 to Table 12. Specifically, in order to more intuitively illustrate the reasoning process and tool usage in agentic RL, we selected one representative case generated by our method from each of the seven benchmarks used in the main experiments above. We use different colored boxes to describe the different states and actions of the agentic LLM during the reasoning process. For multi-hop QA datasets, which mainly include more complex question-answering problems requiring multi-step reasoning, agentic LLMs typically exhibit more multi-turn tool calls and interactions compared to single-hop QA questions. This is consistent with the statistics shown in Figure 4.

Example #1 from HotpotQA Dataset

Question:

What was the father of Kasper Schmeichel voted to be by the IFFHS in 1992?

Labeled Answer:

World's Best Goalkeeper

Model Output by AT²PO(Qwen3-8B)

<think> Okay, let's try to figure out what the father of Kasper Schmeichel was voted to be by the IFFHS in 1992. First, I need to identify who Kasper Schmeichel's father is. I know that Kasper Schmeichel is a Danish footballer, a goalkeeper. So, maybe his father is also a footballer? I should start by searching for Kasper Schmeichel's father. Let me use the Wikipedia search tool for that. Let me search "Kasper Schmeichel father". **</think>**

<search> Kasper Schmeichel father **</search>**

<result>

Page 1: "Kasper Schmeichel" surpassing his father's record. During the Round of 16 match against Croatia, Schmeichel was named Man of the Match after saving a penalty-kick taken by Luka Modrić in extra-time and two penalties during the subsequent shootout, although his opposing counterpart Danijel Subašić saved three penalties in that shootout to eliminate Denmark. Schmeichel is the son of former Denmark international goalkeeper Peter Schmeichel. Having spent most of his childhood living in England due to his father's playing career, Kasper is a native speaker of both Danish and English, the latter of which he speaks with a Northern-English accent.

Page 2: "Peter Schmeichel" "125 greatest living footballers", at the FIFA 100 celebrations. His son, Kasper, is also a professional football goalkeeper, currently playing for Premier League side Leicester City and the Danish national team. Peter Schmeichel was born in the Søborggård parish of Gladsaxe, Denmark, to Inger, a Danish nurse, and Antoni Schmeichel, a Polish jazz musician. He held Polish citizenship until 1970 when he, his father, and his three sisters became Danish citizens. Schmeichel inherited his middle name – Bolesław – from his great-grandfather. He spent his early years in the town of Buddinge, Copenhagen.

Page 3: "Kasper Schmeichel" Hall Grammar School in Cheadle Hulme, Stockport, Greater Manchester, and he also attended Saint Julian's School in Portugal while his father played for Sporting CP in Lisbon. When Alex Bruce, son of Kasper's father's former Manchester United teammate Steve Bruce, joined Kasper at Leeds United, the move saw the two players re-unite after both had previously played football together on the streets as children after school. Schmeichel has two children: a son, Max (born 2010) and a daughter, Isabella (born 2012). Schmeichel signed an ambassadorial deal with Precision Goalkeeping on 17 August 2010.

</result>

<think> Okay, so the search result says that Kasper Schmeichel's father is Peter Schmeichel. Now, the question is about what Peter Schmeichel was voted to be by the IFFHS in 1992. I need to check the IFFHS awards or something related to Peter Schmeichel in 1992. First, I should search who Peter Schmeichel is. Wait, I already have that from the result. Let me search for IFFHS and Peter Schmeichel in 1992. **</think>**

<search> Peter Schmeichel IFFHS 1992 **</search>**

<result>

One turn search results omitted here

<result>

<think> Okay, so the search result says that Peter Schmeichel was voted the IFFHS World's Best Goalkeeper in 1992. The father of Kasper Schmeichel is Peter Schmeichel, a renowned Danish former professional footballer who played as a goalkeeper. According to the information retrieved, Peter Schmeichel was voted the IFFHS World's Best Goalkeeper in 1992. So the answer is that his father was voted as the IFFHS World's Best Goalkeeper in 1992. **</think>**

<answer>

World's Best Goalkeeper

</answer>

Table 6: An example from AT²PO on HotpotQA dataset, with special symbols used in think content, search queries, returned results and final answer highlighted with **blue box**, **red box**, **green box** and **purple box**, respectively.

Example #1 from Musique Dataset

Question:

Which company owns the manufacturer of Learjet 60?

Labeled Answer:

Bombardier Inc/Bombardier

Model Output by AT²PO(Qwen3-8B)

<think> Okay, so I need to find out which company owns the manufacturer of the Learjet 60. Let me start by recalling what I know. Learjet is a brand of business jets, right? I think Learjet was originally a company that made these jets. But I'm not sure about the current ownership. First, I should figure out who manufactures the Learjet 60. Maybe I should search for the manufacturer of the Learjet 60. Let me do that. I'll use the Wikipedia search tool to look up the Learjet 60. **</think>**

<search> Learjet 60 manufacturer **</search>**

<result>

Page 1: "Learjet 60" Learjet 60 The Learjet 60 is a mid-size cabin, medium-range business jet aircraft manufactured by Bombardier Aerospace in Wichita, Kansas. Powered by two Pratt & Whitney Canada PW305A engines, it has a range (with 4 passengers and 2 crew) of with NBAA reserves, ISA. In July 2012 Bombardier Aerospace announced a temporary "production pause" of the latest variant Learjet 60XR to begin in the fourth quarter of 2012. The Learjet 60 is an improved version of the Learjet 55, with a longer fuselage and more powerful turbofan engines. It first flew on 10 October 1990 and received FAA certification.

Page 2: "Learjet 60" and fractional jet operators. Learjet 60 The Learjet 60 is a mid-size cabin, medium-range business jet aircraft manufactured by Bombardier Aerospace in Wichita, Kansas. Powered by two Pratt & Whitney Canada PW305A engines, it has a range (with 4 passengers and 2 crew) of with NBAA reserves, ISA. In July 2012 Bombardier Aerospace announced a temporary "production pause" of the latest variant Learjet 60XR to begin in the fourth quarter of 2012. The Learjet 60 is an improved version of the Learjet 55, with a longer fuselage and more powerful turbofan engines. It first flew on 10 October 1990.

Page 3: "Learjet 60" It also distinguished as the last legacy Learjet, using a variation of the wing that designer Bill Lear adapted from the Swiss military aircraft, the FFA P-16. The next-generation Learjet was to be the Learjet 85 and was an all-new design by Bombardier Aerospace slated for delivery in 2013. The Learjet 60, while a tremendous performer, also maintains the highest incident/accident rate in its class with most accidents occurring during landing. The use of the original Lear 23 gear and wheels (albeit it with an added brake rotor, bringing the total to 3 per wheel assembly).

</result>

<think> Okay, so the Learjet 60 is manufactured by Bombardier Aerospace. Now the question is asking which company owns the manufacturer. So Bombardier Aerospace is the manufacturer. Therefore, I need to find out who owns Bombardier Aerospace. Let me search for Bombardier Aerospace's parent company. **</think>**

<search> Bombardier Aerospace parent company **</search>**

<result>

One turn search results omitted here

<result>

<think> Okay, according to this, Bombardier Aerospace is a division of Bombardier Inc. So the manufacturer of the Learjet 60 is Bombardier Aerospace, which is owned by Bombardier Inc. Therefore, the answer should be Bombardier Inc. **</think>**

<answer>

Bombardier Inc

</answer>

Table 7: An example from AT²PO on Musique dataset, with special symbols used in think content, search queries, returned results and final answer highlighted with **blue box**, **red box**, **green box** and **purple box**, respectively.

Example #1 from 2WikiMultiHopQA Dataset

Question:

Which film was released earlier, August Underground's Penance or Vasantha Raagam?

Labeled Answer:

Vasantha Raagam

Model Output by AT²PO(Qwen3-8B)

<think> Okay, I need to find out which film was released earlier between "August Underground's Penance" and "Vasantha Raagam." Let me start by recalling what I know about these two films. First, "Vasantha Raagam" might be a South Indian film, possibly from Tamil or Telugu cinema, given the name. I remember that "Vasantha Raagam" is a classic film, maybe from the 1960s or 1970s. But I'm not entirely sure. On the other hand, "August Underground's Penance" sounds like it could be an independent or arthouse film. The title seems a bit more modern. However, I'm not certain about its release date. To get accurate information, I should use the Wikipedia search tool to find the release dates of both films. Let me start by searching for "Vasantha Raagam" to find its release year. **</think>**

<search> Vasantha Raagam **</search>**

<result>

Page 1: "Vasantam" Vasantam is a Tollywood film produced by N.V.Prasad & Sanam Naga Ashok kumar on Sri Sai Deva Productions banner, directed by Vikraman. Starring Venkatesh, Arti Agarwal, Kalyani in lead roles and music is composed by S. A. Rajkumar. The film won for two Nandi Awards. The film opened simultaneously alongside the Tamil version of the film, "Priyamaana Thozhi". Ashok (Venkatesh) and Julie (Kalyani) are childhood friends. Ashok is like a family member of Julie's family. They are very close and believe that friendship is above love, yet they do not have any love feelings towards each other.

Page 2: "Vasantham Vanthachu" Vasantham Vanthachu is a 2007 Tamil romantic comedy film directed by K. Viveka Bharathi. The film features Venkat Prabhu and Nanditha Jennifer in lead roles, with Shanmugasundaram, Vaiyapuri, Ramji, Varalakshmi, Sabitha Anand, Vasu Vikram and Velu Thigal playing supporting roles. The film, produced by Velu Thigal and Amba Umadhevi, had musical score by Kavin Saradha and Raj Shankar, and was released on 14 September 2007. Kumaresan (Venkat Prabhu) is a milk vendor in his village, who hates the marriage and stays away from women. His elderly father (Velu Thigal) wants him to get married as soon as possible.

Page 3: "Vasanth" Vasanth is an Indian film director and screenwriter, working in the Tamil film industry. Following a stint as journalist and after assisting K. Balachander, he made his directorial debut with "Keladi Kanmani" in 1990, which along with the trendsetting "Aasai" (1995), are considered his finest and most successful films. Known for being realistic, his films have been described as a "middle path" between parallel and commercial cinema, revolving around social themes and relationships. Vasanth began his career as a journalist and short story writer, before he began working as an assistant director to K. Balachander on 18 films.

</result>

Four turn search steps omitted here

<think> The search result doesn't mention "Vasantha Raagam" directly. However, given that "August Underground's Penance" is from 2007, and "Vasantha Raagam" might be a 1960s film, I can infer that "Vasantha Raagam" was released earlier. Alternatively, if "Vasantha Raagam" is a 2007 film, then both would be the same. But since the user is asking which was released earlier, and "August Underground's Penance" is 2007, I need to check if "Vasantha Raagam" is older. Since I couldn't find exact info, I'll have to make an educated guess based on the available data. Given that "August Underground's Penance" is from 2007, and "Vasantha Raagam" might be a classic film from the 1960s or 1970s, it's likely that "Vasantha Raagam" was released earlier. Based on the information gathered, "August Underground's Penance" was released in 2007. "Vasantha Raagam" is likely a classic film from the 1960s or 1970s, which would make it older than "August Underground's Penance".

</think>

<answer>

Vasantha Raagam

</answer>

Table 8: An example from AT²PO on 2WikiMultiHopQA dataset, with special symbols used in think content, search queries, returned results and final answer highlighted with **blue box**, **red box**, **green box** and **purple box**, respectively.

Example #1 from Bamboogle Dataset

Question:

Who was the first African American mayor of the most populous city in the United States?

Labeled Answer:

David Dinkins

Model Output by AT²PO(Qwen3-8B)

<think> Okay, so I need to find out who was the first African American mayor of the most populous city in the United States. Let me break this down step by step. First, the most populous city in the United States. I remember that this is usually a city like New York City, but I should confirm that. Let me do a quick search to make sure. **</think>**

<search> most populous city in the United States **</search>**

<result>

Page 1: "New York City" New York City, often called New York City (NYC) or simply New York (NY), is the most populous city in the United States. With an estimated 2017 population of 8,622,698 distributed over a land area of about, New York City is also the most densely populated major city in the United States. Located at the southern tip of the state of New York, the city is the center of the New York metropolitan area, the largest metropolitan area in the world by urban landmass and one of the world's most populous megacities.

Page 2: "Demography of the United States" its southern coast - with particular emphasis on the city of Anchorage - and Hawaii's is centered on the island of Oahu. California and Texas are the most populous states, as the mean center of U.S. population has consistently shifted westward and southward. New York City is the most populous city in the United States. The United States Census Bureau shows a population increase of 0.75% for the twelve-month period ending in July 2012. Though high by industrialized country standards, this is below the world average annual rate of 1.1%.

Page 3: *****more search results omitted here*****

</result>

*****Five turn search steps omitted here*****

<think> Okay, so the most populous city is New York City. Now, the next part is to find the first African American mayor of that city. Let me search for that. **</think>**

<search> first African American mayor of New York City **</search>**

<result>

*****One turn search results omitted here*****

</result>

*****Five turn search steps omitted here*****

<think> Okay, the search result says that David Dinkins was the first and only African American mayor of New York City. The most populous city in the United States is New York City. After searching for the first African American mayor of New York City, the information consistently points to David Dinkins as the first and only African American mayor of the city, serving from 1990 to 1993. So the answer should be David Dinkins.

</think>

<answer>

David Dinkins

</answer>

Table 9: An example from AT²PO on Bamboogle dataset, with special symbols used in think content, search queries, returned results and final answer highlighted with **blue box**, **red box**, **green box** and **purple box**, respectively.

Example #1 from NQ Dataset

Question:

how many episodes are there in dragon ball z?

Labeled Answer:291 episodes/291

Model Output by AT²PO(Qwen3-8B)

<think> Okay, so the user is asking how many episodes there are in Dragon Ball Z. Hmm, I need to figure this out. Let me start by recalling what I know. Dragon Ball Z is a popular anime series that's part of the Dragon Ball franchise. It's known for having a lot of episodes, maybe over 200? But I need to confirm that. Let me use the Wikipedia search tool to get accurate data. First, I'll search for "how many episodes are there in dragon ball z" using the tool. **</think>**

<search> how many episodes are there in dragon ball z **</search>**

<result>

Page 1: "Dragon Ball Z" It premiered in Japan on Fuji Television on April 26, 1989, taking over its predecessor's time slot, and ran for 291 episodes until its conclusion on January 31, 1996. Because Toriyama was writing the manga during the production of the anime, "Dragon Ball Z" added original material not adapted from the manga, including lengthening scenes or adding new ones, and adding new attacks and characters not present in the manga. For example, Toriyama was asked to create an additional character for Goku's training with King Kai, resulting in the cricket Gregory.

Page 2: "Dragon Ball" April 5, 2009, the series premiered in Japan airing in Fuji TV. "Dragon Ball Z Kai" reduced the episode count to 159 episodes (167 episodes internationally), from the original footage of 291 episodes. Damaged frames were removed, resulting in some minor shots being remade from scratch in order to fix cropping, and others to address continuity issues. The majority of the international versions, including Funimation Entertainment's English dub, are titled "Dragon Ball Z Kai". premiered on Fuji TV on February 7, 1996 and ran until November 19, 1997 for 64 episodes.

Page 3: "Dragon Ball" both based on the same outline by Toriyama for the story of "Super". Toei Animation produced an anime television series based on the first 194 manga chapters, also titled "Dragon Ball". The series premiered in Japan on Fuji Television on February 26, 1986 and ran until April 19, 1989, lasting 153 episodes. Instead of continuing the anime as "Dragon Ball", Toei Animation decided to carry on with their adaptation under a new name and asked Akira Toriyama to come up with the title.

</result>

<think> Okay, the first result says that Dragon Ball Z ran for 291 episodes from April 26, 1989, to January 31, 1996. That seems like the main series. But the second result mentions "Dragon Ball Z Kai" which reduced the episode count to 159 episodes (167 internationally). However, the user is asking about the original Dragon Ball Z, not the Kai version. The third result talks about the original Dragon Ball series, which had 153 episodes, but that's before Z. So the answer is 291 episodes.

</think>

<answer>

291 episodes

</answer>

Table 10: An example from AT²PO on NQ dataset, with special symbols used in think content, search queries, returned results and final answer highlighted with **blue box**, **red box**, **green box** and **purple box**, respectively.

Example #1 from TriviaQA Dataset

Question:

Who had an 80s No 1 hit with Hold On To The Nights?

Labeled Answer:

Richard Noel Marx/Richard Marx

Model Output by AT²PO(Qwen3-8B)

<think> Okay, so the question is asking who had an 80s No. 1 hit with "Hold On To The Nights." Hmm, I need to figure this out. Let me start by recalling if I know the answer. I think it might be a band from the 80s, maybe a British group? Or maybe a singer? Let me check using the Wikipedia search tool to confirm. First, I'll search for "Hold On To The Nights 80s No 1 hit" to see if that gives me any clues. **</think>**

<search> who had an 80s no 1 hit with hold on to the nights **</search>**

<result>

Page 1: "Hold On to the Nights" Hold On to the Nights "Hold On to the Nights" is a power ballad by American rock singer/songwriter/musician Richard Marx. This was the fourth and final single released from his self-titled debut album, and the first to reach number-one on the US "Billboard" Hot 100 chart. The song has been re-released on numerous albums and is included on Marx's live performance DVD "A Night Out with Friends" (2012). Written by Richard Marx, "Hold On to the Nights" reached the "Billboard" Hot 100 number 1 position on July 23, 1988, preventing Def Leppard's "Pour Some Sugar on Me" from reaching the

Page 2: "Richard Marx (album)" Richard Marx (album) Richard Marx is the debut studio album by singer/songwriter and record producer/arranger, Richard Marx, released in June 1987. Richard became the first male solo artist (and second solo artist overall - the first being Whitney Houston) in recording history to reach the top three of the "Billboard" Hot 100 with four singles from a debut album, including his first number one Pop single "Hold On to the Nights" in mid-1988. The album's first single, "Don't Mean Nothing", hit #1 on "Billboard"'s Hot Mainstream Rock Tracks chart on July 4, 1987.

Page 3: *****more search results omitted here*****

</result>

<think> Okay, according to the search result, "Hold On to the Nights" is a power ballad by Richard Marx, and it reached number one on the US Billboard Hot 100 chart on July 23, 1988. The result also mentions that it was the fourth and final single from his self-titled debut album. So the answer should be Richard Marx. Let me confirm the time frame. The result says it was in mid-1988, which is within the 80s. So the answer is Richard Marx. But let me check if there's any other artist mentioned in the search results. The third result is about the "Family Guy" episode, but that's not relevant. The second result is about the album, which also states that the first number one was "Hold On to the Nights" in mid-1988. So yes, Richard Marx is the correct answer.

</think>**<answer>**

Richard Marx

</answer>

Table 11: An example from AT²PO on TriviaQA dataset, with special symbols used in think content, search queries, returned results and final answer highlighted with **blue box**, **red box**, **green box** and **purple box**, respectively.

Example #1 from PopQA Dataset

Question:

What is Eleanor Davis's occupation?

Labeled Answer:cartoonist/graphic artist/ animator/illustrator

Model Output by AT²PO(Qwen3-8B)

<think> Okay, so I need to find out what Eleanor Davis's occupation is. Let me start by recalling if I know anything about her. Maybe she's an author, artist, or something else? Hmm, not sure. Let me use the Wikipedia search tool to look her up.

</think>

<search> What is Eleanor Davis's occupation? **</search>**

<result>

Page 1: "Eleanor Layfield Davis" Eleanor Layfield Davis Eleanor Layfield Davis (1911–1985), also called ELDA, was an American painter. She served on the Board of Trustees for Meredith College and both Meredith and Wake Forest University award art scholarships in her memory. Eleanor Layfield Davis was born Eleanor Layfield in 1911, in Richmond, Virginia. She grew up in the Raleigh, North Carolina area, attending Peace College and Sweet Briar College. Davis went on to graduate from Meredith College in 1932. She was married to Egbert Lawrence Davis, Jr. (1911–2006) and had four children. Davis did not start painting until she was 47 years old.

Page 2: "Eleanor Layfield Davis" and in 15 individual exhibitions around the Southeast. Posthumous exhibitions include a retrospective at Wake Forest University's Scales Fine Arts Center Gallery in 1986, and "ELDA – Paintings by Eleanor Layfield Davis" at the Sawtooth Center for Visual Arts in 2012. In early 2010, some of her works were included in a Meredith College alumnae art showcase. Some of Davis artwork was purchased by the North Carolina Art Society for the loan collection of the North Carolina Museum of Art. Eleanor Layfield Davis Eleanor Layfield Davis (1911–1985), also called ELDA, was an American painter.

Page 3: "Eleanor Layfield Davis" primarily signed her works with the name ELDA developed from her initials. In addition to paintings, Davis also drew sketches and sculpted. She spent much of her life helping build up the arts community in Winston-Salem, North Carolina. She studied art at the Arts and Crafts Association of Winston-Salem as well as taking some courses at the University of North Carolina at Greensboro. Davis was a founder and President of Associated Artists of Winston-Salem as well as organizer of Art Gallery Originals. Davis died in 1985. During her lifetime, Davis produced many works that she exhibited in group shows.

</result>

Two turn search steps omitted here

<think> Oh, there's another result about Eleanor Davis, born in 1983, an American cartoonist and illustrator. So there are two people named Eleanor Davis. The first is Eleanor Layfield Davis, a painter, and the second is Eleanor Davis, a cartoonist. The user might be referring to the cartoonist, given that the search result for "Eleanor Davis" includes that. Let me check which one is more likely. Since the user didn't specify, but the first result is about the painter, but the second is about a cartoonist. Maybe the answer depends on which one is being asked. However, the initial search query was "Eleanor Davis's occupation" and the result includes both. But the second result is more recent, so maybe that's the one. Let me confirm. The second result says Eleanor Davis (born 1983) is an American cartoonist and illustrator. So if the user is referring to that person, the answer is cartoonist and illustrator. But the first result is about a painter. Maybe the user is referring to the cartoonist, as that is more common in modern contexts. However, without more context, it's ambiguous. But since the search result includes both, perhaps the answer is painter. However, the second result is more recent and has more specific info. Let me check the exact name. The result says "Eleanor Davis (born January 16, 1983) is an American cartoonist and illustrator..." So if the user is referring to that person, the answer is cartoonist and illustrator. But the initial search result is about the painter. Maybe the user is referring to the cartoonist. However, the first result is about Eleanor Layfield Davis, which is a different name. Maybe the user made a typo or is referring to the cartoonist. Since the second result is more recent, I think that's the intended person. Let me confirm. Alternatively, maybe the user is referring to the painter. But since the search result includes both, I need to check which one is more likely. Since the user didn't specify, but the first result is about a painter, and the second is about a cartoonist, perhaps the answer is painter. However, the second result is more detailed. Maybe the answer is cartoonist. But without more info, I need to choose based on the search results. Since the second result is from the search, and it's more recent, I think that's the one. So the answer is cartoonist and illustrator.

</think>

<answer>

Richard Marx

</answer>

Table 12: An example from AT²PO on PopQA dataset, with special symbols used in think content, search queries, returned results and final answer highlighted with **blue box**, **red box**, **green box** and **purple box**, respectively.