

Exploring Reasoning Reward Model for Agents

Kaixuan Fan^{1,2} Kaituo Feng^{1,2} Manyuan Zhang^{2*} Tianshuo Peng¹ Zhixun Li³
Yilei Jiang^{1,2} Shawn Chen² Peng Pei² Xunliang Cai² Xiangyu Yue^{1†}

¹MMLab, CUHK ²Meituan ³SEEM, CUHK

Repository: <https://github.com/kxfan2002/Reagent>

Abstract

Agentic Reinforcement Learning (Agentic RL) has achieved notable success in enabling agents to perform complex reasoning and tool use. However, most methods still relies on sparse outcome-based reward for training. Such feedback fails to differentiate intermediate reasoning quality, leading to suboptimal training results. In this paper, we introduce **Agent Reasoning Reward Model (Agent-RRM)**, a multi-faceted reward model that produces structured feedback for agentic trajectories, including (1) an explicit reasoning trace, (2) a focused critique that provides refinement guidance by highlighting reasoning flaws, and (3) an overall score that evaluates process performance. Leveraging these signals, we systematically investigate three integration strategies: **Reagent-C** (text-augmented refinement), **Reagent-R** (reward-augmented guidance), and **Reagent-U** (unified feedback integration). Extensive evaluations across 12 diverse benchmarks demonstrate that Reagent-U yields substantial performance leaps, achieving 43.7% on GAIA and 46.2% on WebWalkerQA, validating the effectiveness of our reasoning reward model and training schemes. Code, models, and datasets are all released to facilitate future research.

1 Introduction

Reinforcement Learning with Verifiable Reward (RLVR) has achieved remarkable success in improving the reasoning capabilities of Large Language Models (LLMs) (Liu et al., 2025b; Feng et al., 2025b; Tang et al., 2025; Chen et al., 2025b,c). Motivated by this progress, recent works have extended this paradigm to agents, demonstrating its potential to handle complex interactions with dynamic environments and external knowledge sources (Jin et al., 2025; Wu et al., 2025a; Li et al., 2025b).

However, previous agentic RL methods typically rely on sparse, outcome-based rewards based solely on final correctness (Jin et al., 2025; Wu et al., 2025a; Li et al., 2025b). This design is inherently limiting for long-horizon agentic tasks requiring multi-step tool utilization (Feng et al., 2025c; Liu et al., 2025a; Zhang et al., 2025e). In such settings, outcome-based supervision fails to differentiate high-quality intermediate reasoning from entirely incorrect attempts, for instance, treating a trajectory that fails only at the final step as a total failure. Consequently this coarse-grained binary supervision obscures the value of successful intermediate steps, resulting in sub-optimal performance (Dong et al., 2025).

To provide more granular feedback, recent research has pivoted toward integrating Reward Models into Agentic RL. However, the effective deployment of Reward Models remains hampered by two bottlenecks. First, while step-level rewards offer finer granularity feedback (Xi et al., 2025; Liu et al., 2025a; Xu et al., 2025), they are often plagued by prohibitive annotation costs (Rahman et al., 2025) and a susceptibility to reward hacking (Zhang et al., 2025c). Second, existing reasoning-based Reward Models focus on pair-wise preferences (Li et al., 2025c; Liu et al., 2025a; Hu et al., 2025), which frequently introduces inherent biases and fails to capture fine-grained quality gradations between trajectories or provide actionable guidance for refinement (Jian et al., 2025; Zhang et al., 2025d). Furthermore, most of these efforts exclusively rely on numeric reward feedback for training, leaving the natural language critique (Zhang et al., 2025a) largely unexplored, which could provide more granular guidance for agentic policy.

To this end, we develop **Agent Reasoning Reward Model (Agent-RRM)**, a multi-faceted evaluator designed to provide reasoning-aware feedback for agentic trajectories. Unlike conventional Reward Models that yield merely scalar scores or

*Project Leader.

†Corresponding Author.

binary preferences, Agent-RRM conducts explicit reasoning to justify its assessments. For each trajectory, it generates a decomposed judgment comprising: (1) an internal reasoning trace that analyzes logical consistency of trajectory; (2) a targeted critique identifying specific flaws to guide refinement; and (3) a holistic quality score. This hierarchy of signals provides dense, multi-dimensional supervision, combining scalar rewards for global optimization with textual critiques for explicit error correction—all without necessitating ground truth.

Building upon these informative signals, we perform a systematic investigation into the integration of Agent-RRM and Agentic RL. We formalize this integration through a unified scheme with three variants: Text-augmented Refinement, where agents polish trajectories based on Agent-RRM’s textual feedback; Reward-augmented Guidance, which complements rule-based rewards with model-based signals; and Unified Feedback Integration, which harmonizes multi-source rewards with critique-augmented sampling. We denote the agent policy models of these variants as **Reagent-C**, **Reagent-R**, and **Reagent-U**, respectively. Notably, our experiments demonstrate that Reagent-U achieves superior performance by synthesizing these feedback modalities, reaching 43.7% on GAIA and 46.2% on WebWalkerQA. Our study provides a comprehensive roadmap for harnessing multi-level feedback to accelerate agentic RL.

To support this investigation, we curate **four specialized datasets** that provide high-quality trajectories for both agent reasoning and reward model training. Extensive experiments across **12 diverse benchmarks** demonstrate that Reagent models achieve significant performance gains, underscoring the efficacy of multi-level reasoning-based feedback signals in complex agentic tasks.

In summary, our contributions are as follows:

- We introduce **Agent-RRM**, a multi-faceted evaluator that generates structured feedback including explicit reasoning rationales, actionable critiques, and holistic quality scores, providing a transparent and granular assessment.
- We systematically explore three agent variants with Agent-RRM: Text-augmented Refinement (**Reagent-C**), Reward-augmented Guidance (**Reagent-R**), and Integrated Feedback Optimization (**Reagent-U**). This provides a roadmap for using reasoning rewards to enhance agent performance.

- We curate and release **four high-quality datasets** specifically tailored for training reasoning agent and reward model. These resources provide the community with valuable assets to advance research in multi-granular feedback for agentic reinforcement learning.

2 Related Work

2.1 Agentic Reinforcement Learning

Agentic Reinforcement Learning (Agentic RL) has emerged as a cornerstone for developing agents capable of operating in dynamic, open-ended environments (Dong et al., 2025; Lù et al., 2025). Recent advancements (Jin et al., 2025; Wang et al., 2025; Xia et al., 2025; Wu et al., 2025c; Li et al., 2025f; Song et al., 2025) illustrate that RL can effectively instill multi-step information-seeking and tool-use proficiencies. For example, Search-R1 (Jin et al., 2025) demonstrates that agentic RL enables LLMs to interleave multi-turn web search, substantially improving retrieval-augmented reasoning performance. WebSailor (Li et al., 2025b) further shows that agentic RL can scale to long-horizon web navigation, equipping agents with the ability to reduce extreme uncertainty in complex information-seeking tasks. Agent0 (Xia et al., 2025) introduces a co-evolutionary process where tool-aware reasoning behaviors emerge without human-curated supervision. Despite these successes, most existing methods rely heavily on sparse, outcome-based rewards, which often limits training efficacy and hampers agent’s ability to rectify intricate intermediate errors (Dong et al., 2025; Lin et al., 2025).

2.2 Reward Modeling

Reward Models (RMs) play a central role in optimizing complex reasoning tasks by providing learning signals for policy improvement (Wang et al., 2024; Fan et al., 2025; Li et al., 2025a). Motivated by Deepseek-R1 (Guo et al., 2025), reasoning-aware reward models are introduced to perform explicit reasoning before reward assignment for delivering higher-quality and more transparent supervision (Whitehouse et al., 2025; Zhang et al., 2025b). For instance, RM-R1 (Chen et al., 2025d) introduces a generative reasoning-based reward model to first derive explicit reasoning rubrics and then evaluate candidate responses accordingly. R1-Reward (Zhang et al., 2025b) proposes a multi-modal reasoning reward model and introduces a stabilized RL algorithm that improves training ro-

SFT Data	DeepScaler: 8k	LIMO: 0.6k	OpenMathReasoning: 2k	SimpleRL_Zoo: 2k
APPS: 3k	Tool_Star_Math: 1k	TaskCraft: 0.9k	Tool_Star_QA: 16k	SimpleDeepSearcher: 0.4k
Sky_T1: 4k	Heysquad: 2k	LongAudio: 5k	OCR-VQA: 5k	PixelReasoner: 1k
WebWalkerQA: 0.4k	file_Heysquad: 1k	file_OCR_VQA: 2k	file_zip: 0.4k	
Tool-VQA: 1.0k				
RL Data	DeepScaleR: 19k	Deepmath: 96k	LIMO: 0.8k	OpenMathReasoning: 7k
APPS: 5k	SimpleRL_Zoo: 17k	Tool_star_rl: 5k	Hiersearch: 28k	ASearcher-Base-35k: 35k
MetaMathQA: 395k	Sailorfog-QA: 0.02k	Sky-T1: 10k	Tool_star_rl: 5k	SimpleDeepSearcher: 0.9k
Retool_sft: 2k	Webagent_sft: 5k	WebShaper: 0.5k	WebWalkerQA: 2k	HeySQuAD: 4k
TaskCraft: 13k	LiveVQA: 8k	MMK12: 3k	OCR_VQA: 8k	PixelReasoner: 2k
LongAudio: 10k	ToolVQA: 9k	file_LiveVQA: 8k	file_LongAudio: 10k	
ScienceQA: 1k				

Figure 1: Detailed distribution information of Reagent-SFT-55.6K and Reagent-RL-709K.

bustness for multimodal RMs. In agent domain, reasoning-based reward models still remain under-explored. Atom-Searcher (Deng et al., 2025) directly utilizes a Qwen3-30B-A3B without training as reward model to assign scores to agent steps, while PPR (Xu et al., 2025) employs a process reward model to evaluate trajectory steps based on a predefined principle set. However, these methods remain confined to step-level scalar rewards, which are susceptible to reward hacking and fail to provide language-based guidance necessary for rectifying complex logic flaws.

3 Method

3.1 Preliminaries: GRPO Framework

In Group Relative Policy Optimization (GRPO) (Shao et al., 2024), for a query q sampled from the dataset $P(q)$, the policy π_θ generates a group of G outputs $\{o_i\}_{i=1}^G$ such that:

$$o_i \sim \pi_{\theta_{old}}(o|q). \quad (1)$$

Let $r_i(\theta) = \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}$ denote the importance sampling ratio. The GRPO objective is formulated as:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{q \sim P(q)} \{ \mathbb{E}_{\{o_i\} \sim \pi_{\theta_{old}}} \left[\frac{1}{G} \sum_{i=1}^G \left(\min(r_i(\theta)A_i, \text{clip}_\epsilon) - \beta \mathbb{D}_{KL}^{(i)} \right) \right], \quad (2)$$

where clip_ϵ denotes $\text{clip}(r_i(\theta), 1 - \epsilon, 1 + \epsilon)A_i$, and $\mathbb{D}_{KL}^{(i)}$ denotes the KL divergence between current policy and reference model π_{ref} for the i -th output:

$$\mathbb{D}_{KL}^{(i)} = \mathbb{D}_{KL}(\pi_\theta(o_i|q) || \pi_{ref}(o_i|q)). \quad (3)$$

The advantage A_i is computed by normalizing the rewards within the group $\mathbf{R} = \{R_1, \dots, R_G\}$:

$$A_i = \frac{R_i - \text{mean}(\mathbf{R})}{\text{std}(\mathbf{R})}. \quad (4)$$

where R_i is the reward for the i -th output assigned by the reward system.

3.2 Agentic Tool Design

To enable effective interaction with diverse environments, we equip the agent with a suite of six specialized tools covering information retrieval, code execution, and multi-modal perception:

- **Search:** Given a query, retrieve relevant search results using the Bing search engine.
- **Web Browse:** Given a URL and a query, fetch the webpage content and generate a response to the query based on the page information.
- **Python Code Interpreter:** Execute a provided Python code snippet and return the execution results.
- **File Reader:** Access and extract a file and return its textual content.
- **Image Descriptor:** Given an image and a query, generate a textual response to the query conditioned on the visual features.
- **Audio Converter:** Transcribe an input audio file into text.

3.3 Dataset Construction

We curate four specialized datasets spanning mathematical deduction, multimodal understanding, web-based information seeking, and complex tool utilization. These datasets support the distinct training requirements of both Reagent and Agent-RRM.

Agent Training Datasets We synthesize diverse QA benchmarks to enhance the agent’s reasoning and tool-use capabilities. To ensure data quality, we apply a rigorous three-stage pipeline: (1) filtering samples with ambiguous ground truths; (2) cross source deduplication; and (3) difficulty-aware sam-

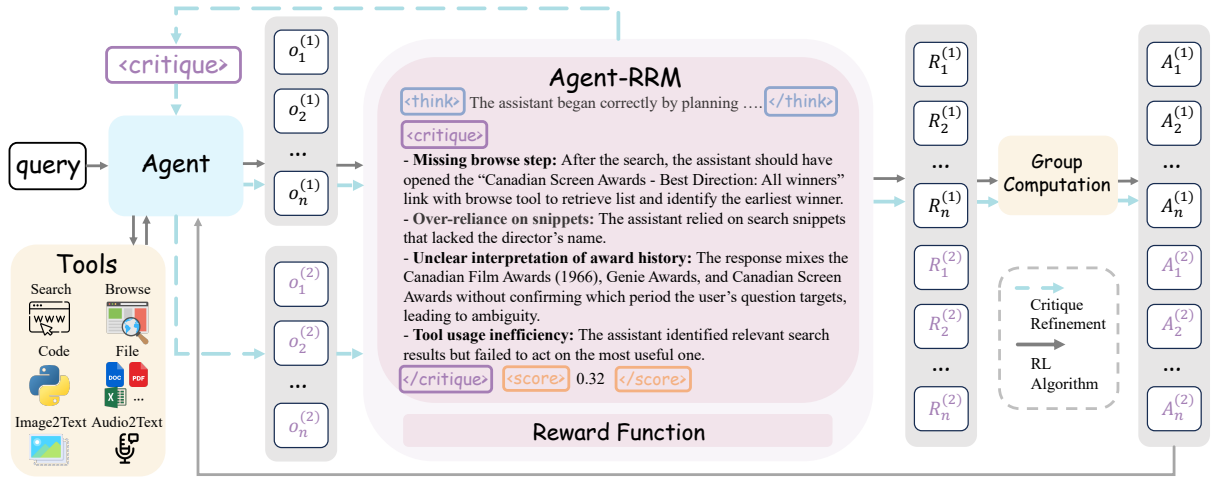


Figure 2: Overview of the **Reagent** training scheme. We explore three integration variants: **Reagent-C** (blue arrows), **Reagent-R** (gray arrows), and **Reagent-U** (both arrows).

pling. This yields **Reagent-RL-709K**, a comprehensive corpus of 709k question-answer pairs for RL training. For Supervised Fine-tuning (SFT), we prioritize the holistic quality of reasoning trajectories. Using DeepSeek-V3.1, we generate and retain only trajectories that lead to correct final answers, resulting in **Reagent-SFT-55.6K** as high-quality cold-start data. Dataset distribution is shown in Figure 1. Specific selection criteria and filtering thresholds are detailed in Appendix A.2.

Reward Model Datasets Training a robust Agent-RRM necessitates exposure to a wide spectrum of logical error patterns and diverse response styles. Building upon the Reagent-RL-709K, we construct two meticulously labeled datasets: **Reagent-RRM-SFT-28K** and **Reagent-RRM-RL-90K** for SFT and RL stage respectively. We sample reasoning trajectories from an ensemble of models including Qwen3-8B/14B, Qwen3-ARPO-DeepSearch (8B/14B), Qwen2.5-7B-ARPO, Qwen2.5-WebDancer (7B/32B), and DeepSeekV3.1 to maximize the coverage of potential error patterns. These trajectories are then annotated by GPT-OSS-120B to generate structured three-part judgments: an analytical <think> trace, a targeted <critique> of flaws, and a holistic <score> ($s \in [0, 1]$). See Appendix A.3 for prompts and process details.

3.4 Agent-RRM: Reward Model Training

To facilitate granular feedback, we train **Agent Reasoning Reward Model (Agent-RRM)** to generate multi-dimensional judgments consisting of three components: <think>, an internal reasoning

trace analyzing trajectory quality, <critique>, a targeted identification of reasoning or execution flaws, and <score>, a scalar quality assessment within $[0, 1]$. Following (Chen et al., 2025d; Zhang et al., 2025b), we adopt a two-stage training procedure. First, we conduct SFT on Reagent-RRM-SFT-28K to instill the structured output format and foundational evaluative capabilities. Subsequently, we apply GRPO on Reagent-RRM-RL-90K to refine the model’s evaluative rationales and ensure the calibration of its scalar rewards. This training paradigm ensures that Agent-RRM can generate high-fidelity, self-consistent feedback even in the absence of ground-truth answers, making it highly effective for complex, open-ended agentic tasks.

3.5 Reagent: Integrating Reasoning Rewards into Agents

In this section, we introduce our agent policy model **Reagent**, and present three variants that explore different ways of incorporating reasoning rewards and critiques into agentic policies.

To provide a robust starting point for RL, we fine-tune the base model on the **Reagent-SFT-55.6K** dataset. This stage ensures the agent acquires fundamental reasoning and tool-calling proficiencies. The resulting optimized policy, $\pi_{\theta_{SFT}}$, serves as the seed model for subsequent RL investigations in Section 3.5.2 and 3.5.3. We investigate three Reagent variants to explore the synergy between Agent-RRM and the agent: (1) **Textual-augmented Refinement (Reagent-C)**, which evaluates the immediate utility of textual critiques via zero-shot, in-context refinement; (2) **Reward-augmented Guidance (Reagent-R)**, which optimizes $\pi_{\theta_{SFT}}$ by com-

plementing rule-based rewards with model-based scalar signals; and (3) **Unified Feedback Integration (Reagent-U)**, which harmonizes both modalities within a joint optimization loop. The overall framework is illustrated in Figure 2.

3.5.1 Textual-augmented Refinement

This variant (**Reagent-C**) exploits textual critiques from Agent-RRM for training-free refinement, applied directly to the Qwen3-8B via in-context prompting.

For each query q , the agent first generates an initial response $o_i^{(1)} \sim \pi_\theta(o|q)$. Subsequently, Agent-RRM analyzes $o_i^{(1)}$ to produce a targeted critique c_i via its `<critique>` component, identifying specific reasoning flaws or execution errors. The agent then performs a refined pass conditioned on feedback:

$$o_i^{(2)} \sim \pi_\theta(o|q, o_i^{(1)}, c_i), \quad (5)$$

where the augmented context $(q, o_i^{(1)}, c_i)$ provides the original task and explicit guidance for correction. Crucially, the policy π_θ remains frozen in this variant, allowing us to isolate and evaluate the agent’s in-context refinement capability. All reported results for Reagent-C correspond to the refined outputs $\{o_i^{(2)}\}$.

3.5.2 Reward-augmented Guidance

This variant (**Reagent-R**) utilizes the scalar score from Agent-RRM to provide fine-grained quality assessments of agent trajectories. Following standard GRPO sampling procedure, the agent generates G outputs $o_i \sim \pi_{\theta_{old}}(o|q)$. The reward R_i is defined as a combination of rule-based correctness and model-based quality evaluation:

$$R_i = R_{rule}(q, o_i) + \lambda \cdot R_{model}(q, o_i), \quad (6)$$

where R_{rule} validates final answer correctness, R_{model} is extracted from Agent-RRM’s `<score>`, and λ is a scaling factor balancing their contributions. This variant alleviates the sparsity of rule-based rewards by providing reasoning-aware feedback. It enables the agent to capture a fine-grained spectrum of trajectory quality, effectively rewarding logical merit while penalizing reasoning deficiencies regardless of final answer’s correctness.

3.5.3 Unified Feedback Integration

This variant (**Reagent-U**) harmonizes scalar rewards and textual critique-driven refinement within a unified RL loop. By simultaneously optimizing

initial generation quality and refinement capability, we investigate whether these objectives can yield synergistic improvements through mutual reinforcement for agent.

For each query q , the agent performs a two-stage sampling:

$$o_i^{(1)} \sim \pi_{\theta_{old}}(o|q), \quad o_i^{(2)} \sim \pi_{\theta_{old}}(o|q, o_i^{(1)}, c_i), \quad (7)$$

where $o_i^{(1)}$ is the initial attempt and $o_i^{(2)}$ is the refined response guided by `<critique>` c_i generated by Agent-RRM.

We pool all trajectories from both stages into $\mathcal{G}_{pool} = \{o_i^{(k)} \mid i \in [G], k \in \{1, 2\}\}$ and compute combined reward $R_i^{(k)}$ via Eq. 6. The advantage is computed across this unified pool:

$$A_i^{(k)} = \frac{R_i^{(k)} - \text{mean}(\mathbf{R}_{pool})}{\text{std}(\mathbf{R}_{pool})}, \quad (8)$$

where $\mathbf{R}_{pool} = \{R_i^{(k)} \mid o_i^{(k)} \in \mathcal{G}_{pool}\}$. The unified objective is formulated as:

$$\mathcal{J}_U(\theta) = \mathbb{E} \left[\frac{1}{2G} \sum_{k=1}^2 \sum_{i=1}^G \left(\min(r_i^{(k)}(\theta) A_i^{(k)}, \text{clip}_\epsilon) - \beta \mathbb{D}_{KL}^{(i,k)} \right) \right], \quad (9)$$

where the importance ratio $r_i^{(k)}(\theta)$ and KL penalty $\mathbb{D}_{KL}^{(i,k)}$ are computed relative to their respective contexts. By normalizing advantages across all initial and refined trajectories, Reagent-U encourages the agent to optimize for overall trajectory quality, effectively boosting the agent’s core reasoning and tool-calling performance. Notably, textual critiques are utilized exclusively during the training phase to internalize reasoning capabilities; at inference time, ReAgent-U operates as a standard agent without additional critique refinement or external guidance.

4 Experiments

Benchmarks We evaluate comprehensively on multiple challenging benchmarks. **(1) Mathematical Reasoning:** AIME24 (Zhang and Math-AI, 2024), AIME25 (Zhang and Math-AI, 2025), GSM8K (Cobbe et al., 2021) and MATH500 (Lightman et al., 2023). **(2) Knowledge-Intensive Reasoning:** HotpotQA (Yang et al., 2018), 2Wiki (Ho et al., 2020), Bamboogle (Press et al., 2023) and MuSiQue (Trivedi et al., 2022). **(3) General Agent and Search Reasoning:** GAIA (Mialon et al., 2023), WebWalkerQA (Wu et al., 2025b), Humanity’s Last Exam (HLE) (Phan et al., 2025) and xbench (Chen et al., 2025a).

Table 1: Comprehensive Evaluation on General Agent and Search Benchmarks.

Method	Backbone	GAIA (text)				WebWalkerQA				HLE	xbench
		Lv.1	Lv.2	Lv.3	Avg.	Easy	Med.	Hard	Avg.	Avg.	Avg.
<i>Proprietary Agents</i>											
-	OpenAI-o3	-	-	-	70.5	-	-	-	71.7	20.2	66.0
-	o1-preview	-	-	-	-	11.9	10.4	7.9	9.9	11.1	-
-	Calude-4-Sonnet	-	-	-	68.3	-	-	-	61.7	20.2	64.0
-	OpenAI DeepResearch	-	-	-	67.4	-	-	-	-	26.6	-
<i>Open-source Baselines ($\leq 8B$)</i>											
WebThinker (Li et al., 2025e)	Qwen3-8B	43.6	11.5	0.0	22.3	6.7	13.1	16.9	13.0	6.6	13.0
WebDancer (Wu et al., 2025a)	Qwen2.5-7B	41.0	30.7	0.0	31.0	40.6	44.1	28.2	36.0	-	-
VerlTool (Jiang et al., 2025)	Qwen3-8B	-	-	-	34.0	-	-	-	-	8.4	-
ARPO (Dong et al., 2025)	Qwen3-8B	53.9	32.7	16.7	38.8	26.7	33.3	29.6	30.5	8.8	25.0
<i>Open-source Baselines ($\leq 32B$)</i>											
-	QwQ-32B	30.9	6.5	5.2	18.9	7.5	2.1	4.2	3.8	6.4	10.0
-	DeepSeek-R1-671B	40.5	21.2	5.2	25.2	5.0	11.8	11.3	10.0	8.6	32.0
Tree-GRPO (Ji et al., 2025)	Qwen2.5-14B	20.8	24.3	7.3	21.0	11.1	15.5	10.8	12.8	-	-
ARPO (Dong et al., 2025)	Qwen3-14B	<u>56.4</u>	<u>40.4</u>	16.7	43.7	31.1	42.9	31.0	36.0	<u>10.0</u>	32.0
Search-o1 (Li et al., 2025d)	QwQ-32B-Preview	53.8	34.6	16.7	39.8	43.1	35.0	27.1	34.1	10.8	40.0
WebDancer (Wu et al., 2025a)	Qwen2.5-32B	46.1	44.2	<u>8.3</u>	<u>40.7</u>	44.3	<u>46.7</u>	29.2	38.4	-	38.0
<i>Open-source Baselines with Process Reward</i>											
Atom-Searcher (Deng et al., 2025)	Qwen2.5-7B	18.0	21.2	0.0	17.5	31.7	23.7	37.0	27.9	<u>10.0</u>	21.0
<i>Our Agents</i>											
-	Qwen3-8B	28.2	21.2	0.0	21.4	31.1	28.6	28.2	29.0	4.0	9.0
Reagent w/o Agent-RRM	Qwen3-8B	41.0	36.5	0.0	34.0	44.4	45.0	41.3	43.5	6.8	32.0
Reagent-C (Direct Inference)	Qwen3-8B	30.8	23.1	16.7	25.2	35.6	38.1	32.4	35.5	4.6	15.0
Reagent-R	Qwen3-8B	51.3	30.8	16.7	36.9	<u>47.5</u>	46.0	<u>42.9</u>	<u>45.3</u>	<u>10.0</u>	<u>41.0</u>
Reagent-U	Qwen3-8B	59.0	38.5	16.7	43.7	49.2	46.8	43.3	46.2	10.8	43.0

Implementation Details Following (Dong et al., 2025; Feng et al., 2025a; Wu et al., 2025a), we employ a two-phase training pipeline: Supervised Fine-Tuning followed by Reinforcement Learning. This training protocol mitigates optimization instability in early RL stages and equips the agent with the foundational skills necessary for effective tool interaction. Both agent models and reward model are initialized from Qwen3-8B (Yang et al., 2025).

Both our agent models and Agent-RRM are trained on 8 NVIDIA A800-80G GPUs. Batch size is set to 32 for both SFT and RL. We use AdamW optimizer. Learning rate is set to 1×10^{-5} for SFT and 5×10^{-7} for RL. λ in Eq. 6 is set to 0.3. Detailed hyperparameters and compute resources are deferred to Appendix B. For evaluation metrics, following (Dong et al., 2025), we utilize Qwen2.5-72B-Instruct as judge model to perform binary scoring based on ground truth and agent prediction. To ensure a rigorous comparison with

prior works (Dong et al., 2025; Wu et al., 2025a), unless otherwise specified, we report pass@1 using a decoding temperature of 0.6 and top-p of 0.95. Evaluation details are shown in Appendix C.

4.1 Can Textual Critiques Guide Inference-Time Refinement?

To investigate the direct impact of textual critiques, we evaluate Reagent-C—a training-free variant—on Qwen3-8B using Agent-RRM for critique guidance. As shown in Table 1 and Table 2, Reagent-C achieves consistent performance gains across all benchmarks without any parameter updates. Improvements are particularly pronounced in Mathematical Reasoning, while solid advancements are also observed in General Agentic and Knowledge-Intensive tasks. We attribute this versatility to Agent-RRM’s diagnostic capacity, which effectively pinpoints logical fallacies and tool-execution errors within complex trajectories. See case study

Table 2: Results on Knowledge-Intensive Reasoning and Math Benchmarks. (HQA: HotpotQA)

Method	Backbone	Knowledge-Intensive Reasoning				Mathematical Reasoning			
		HQA	2Wiki	Bamboogle	MuSiQue	AIME24	AIME25	MATH500	GSM8K
<i>Proprietary Agents</i>									
-	GPT-4o	54.0	49.5	68.8	24.0	13.4	25.7	60.3	-
-	o1-preview	-	-	-	-	46.7	-	85.5	-
-	Claude-4-Sonnet	-	-	-	-	43.4	33.1	93.4	-
<i>Open-source Baselines ($\leq 8B$)</i>									
Search-R1 (Jin et al., 2025)	Qwen2.5-7B	43.3	38.2	43.2	19.6	-	-	-	-
VerlTool (Jiang et al., 2025)	Qwen2.5-7B ¹	42.6	39.2	38.4	18.0	36.7	33.3	82.8	92.1
ARPO (Dong et al., 2025)	Qwen2.5-7B	58.8	76.1	71.5	<u>31.1</u>	33.3	30.0	88.8	92.2
ARPO (Dong et al., 2025)	Qwen3-8B	-	-	-	-	33.3	30.0	88.4	93.4
AgentFlow (Li et al., 2025f)	Qwen2.5-7B	57.0	77.2	69.6	25.3	40.0	-	-	-
<i>Open-source Baselines ($\leq 32B$)</i>									
Search-o1 (Li et al., 2025d)	QwQ-32B-Preview	45.2	58.0	56.0	16.6	<u>56.7</u>	-	86.4	-
Tree-GRPO (Ji et al., 2025)	Qwen2.5-14B	50.2	50.5	54.4	25.9	-	-	-	-
ARPO (Dong et al., 2025)	Qwen3-14B	-	-	-	-	36.7	30.0	83.0	93.6
<i>Open-source Baselines with Process Reward</i>									
Atom-Searcher (Deng et al., 2025)	Qwen2.5-7B	57.3	66.9	70.7	27.6	-	-	-	-
PPR-Instruct (Xu et al., 2025)	Qwen2.5-7B	38.7	31.0	41.2	15.5	-	-	-	-
<i>Our Agents</i>									
-	Qwen3-8B	52.0	58.0	53.6	22.1	46.7	40.0	90.4	94.6
Reagent w/o Agent-RRM	Qwen3-8B	65.8	77.0	61.6	28.1	50.0	43.3	90.8	94.5
Reagent-C (Direct Inference)	Qwen3-8B	61.0	68.9	61.6	25.0	<u>56.7</u>	<u>46.7</u>	93.8	<u>94.9</u>
Reagent-R	Qwen3-8B	<u>67.9</u>	79.0	<u>72.8</u>	28.3	53.3	50.0	<u>92.2</u>	94.1
Reagent-U	Qwen3-8B	68.1	<u>78.8</u>	76.8	31.3	60.0	50.0	93.8	95.1

in Appendix D.

The results confirm that the second response $\{o_i^{(2)}\}$ consistently achieves better performance compared to the initial response $\{o_i^{(1)}\}$. Crucially, the widening margin between the first and second response underscores that many initial failures stem from transient execution errors or logical oversights. Our critiques are uniquely positioned to rectify these flaws by offering precise, actionable feedback. Furthermore, since Agent-RRM operates without access to ground-truth answers, these performance gains empirically validate its capacity to diagnose reasoning flaws and tool-execution errors. This highlights that textual critiques offer the high-granularity supervision essential for mastering complex, multi-step agentic tasks.

4.2 Does Model-based Reward Improve Learning?

To explore whether dense model-based rewards can alleviate reward sparsity in agentic RL, we evaluate

¹VerlTool backbones: Qwen2.5-7B (knowledge) and Qwen2.5-Math-7B (math).

Reagent-R, which augments rule-based outcome rewards with holistic reasoning-level scores from Agent-RRM. As shown in Table 1 and Table 2, Reagent-R consistently outperforms rule-based reward baseline (Reagent w/o Agent-RRM) across all benchmarks. Specifically, Reagent-R achieves 72.8% on Bamboogle and 41.0% on xbench, surpassing Reagent w/o Agent-RRM by 11.2 and 9.0 percentage points, respectively. These results suggest that holistic model-based rewards provide more informative feedback for complex, multi-step reasoning scenarios, where sparse binary outcomes often provide overly coarse and limited guidance for learning.

Reagent-R serves as a critical ablation to isolate the impact of scalar supervision by excluding the textual critiques used in Reagent-U. While Reagent-R consistently outperforms sparse-reward baselines, it remains inferior to Reagent-U across most tasks. This performance gap suggests that while continuous scores better differentiate trajectory quality, they lack the explicit, structural guidance inherent in textual feedback. This indicates a clear need for supervision with richer informational

Table 3: Performance of Reagent-U on GAIA text set and full set (multi-modal).

Model	GAIA (text)		GAIA (full)	
	pass@1	pass@3	pass@1	pass@3
Qwen3-8B (Yang et al., 2025)	21.4	24.3	20.0	26.7
MCP-R1 (Anonymous, 2025)	39.8	52.4	37.6	51.5
Reagent-U	43.7	53.4	38.8	53.9

granularity, such as the textual critiques integrated into Reagent-U.

4.3 Does Unified Feedback Synergistically Boost Performance?

Evaluations in Tables 1 and 2 reveal that the unified feedback mechanism in Reagent-U consistently outperforms all baselines across a diverse spectrum of reasoning and agentic benchmarks. Specifically, Reagent-U achieves 43.7% on GAIA (text) and 46.2% on WebWalkerQA, surpassing all compared methods. Beyond its excellence in general agentic and search tasks, Reagent-U maintains robustness in knowledge-intensive and mathematical reasoning, securing 76.8% on Bamboogle and 60.0% on AIME24. In contrast, many existing baselines fail to generalize across diverse domains, often suffering from significant performance trade-offs. This balanced proficiency indicates that Reagent-U augments multi-tool, multi-turn reasoning capabilities rather than merely optimizing for web search. Such results demonstrate a comprehensive long-horizon decision-making ability, effectively showing that the unified feedback mechanism—integrating both scalar rewards and textual critiques—allows the agent to internalize a more sophisticated policy across complex, heterogeneous tasks.

4.4 Beyond Text-Only: Cross-Modal Reasoning and Complex Tool Use

To evaluate Reagent-U’s proficiency across diverse modalities and tools, we conduct analysis on full GAIA benchmark. While existing studies (Dong et al., 2025; Jiang et al., 2025; Li et al., 2025b,f) focus on the text subset, which emphasizes web navigation and information retrieval ability, we argue that such a narrow scope overlooks the heterogeneous reasoning capabilities required for complex real-world tasks. By evaluating on the GAIA full set, we challenge the agent with tasks requiring the integration of open-domain search, multimodal interpretation, python coding, and file-based rea-

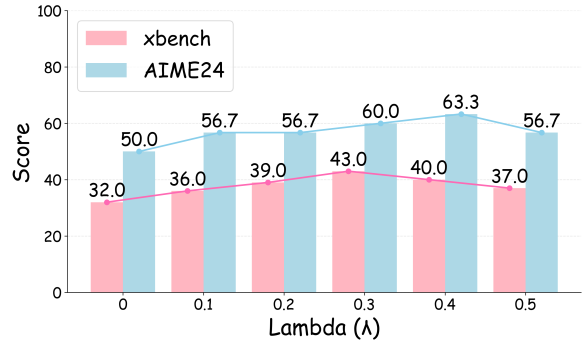


Figure 3: Impact of Agent-RRM reward weight λ on task performance.

soning. As shown in Table 3, Reagent-U not only maintains competitive performance on the text subset but also significantly outperforms baselines on the full set. These results confirm that Reagent-U fosters a versatile agentic intelligence that generalizes across a broad task spectrum rather than overfitting to specific text-based requirements.

4.5 Parameter Analysis on λ

To evaluate the impact of the Agent-RRM reward weight λ , we conduct a parameter analysis on AIME24 (math) and xbench (deep search). Figure 3 shows that agent performance initially increases with rising λ values, demonstrating that the integration of reasoning-based rewards enhances the agent’s decision-making compared to the baseline ($\lambda = 0$). Specifically, performance reaches a plateau between $\lambda \in [0.2, 0.4]$, followed by a slight decline at $\lambda = 0.5$. This trend suggests that while moderate reasoning feedback provides essential supervisory signals, a disproportionately high weight may over-emphasize intermediate steps at the expense of final task completion. Consequently, balancing Agent-RRM rewards with rule-based outcome reward is crucial to maintain an optimal trade-off between reasoning and outcome supervision.

5 Conclusion

In this work, we introduce Agent-RRM, a multi-faceted reasoning reward model designed to provide textual critiques and holistic reasoning-aware reward. Building upon this, we present Reagent, a comprehensive scheme designed to explore the efficacy of multi-dimensional feedback in agentic learning. Our systematic evaluations reveal that while textual critiques provide diagnostic guidance for inference-time refinement, model-based rewards serve to mitigate signal sparsity during training. Together, these signals significantly bolster the

agent’s long-horizon reasoning and multi-step tool-use proficiency, leading to consistent gains across diverse complex, multi-modal tasks.

6 Limitations

We discuss the limitations of our work and potential directions for future research as follows:

First, our current experiments primarily focus on models at the 8B parameter scale. While this setting demonstrates the efficacy of our Reagent scheme, its scaling behavior on larger-scale models remains to be explored. Future work could investigate how more powerful base models might further amplify the benefits of structured reasoning feedback. Second, moving beyond standardized benchmarks to handle broader toolsets and more intricate reasoning chains is essential. Future works can explore open-ended, real-world applications (e.g., AI for science) that involve more diverse toolsets and unpredictable task environments to further validate the scheme’s adaptability.

References

- Anonymous. 2025. [MCP-r1: Generalized real-world task agent mastering dozens of tools](#). In *Submitted to The Fourteenth International Conference on Learning Representations*. Under review.
- Kaiyuan Chen, Yixin Ren, Yang Liu, Xiaobo Hu, Haotang Tian, Tianbao Xie, Fangfu Liu, Haoye Zhang, Hongzhang Liu, Yuan Gong, and 1 others. 2025a. [xbench: Tracking agents productivity scaling with profession-aligned real-world evaluations](#). *arXiv preprint arXiv:2506.13651*.
- Shuang Chen, Yue Guo, Zhaochen Su, Yafu Li, Yulun Wu, Jiacheng Chen, Jiayu Chen, Weijie Wang, Xiaoye Qu, and Yu Cheng. 2025b. [Advancing multimodal reasoning: From optimized cold start to staged reinforcement learning](#). *arXiv preprint arXiv:2506.04207*.
- Shuang Chen, Yue Guo, Yimeng Ye, Shijue Huang, Wenbo Hu, Haoxi Li, Manyuan Zhang, Jiayu Chen, Song Guo, and Nanyun Peng. 2025c. [Ares: Multimodal adaptive reasoning via difficulty-aware token-level entropy shaping](#). *arXiv preprint arXiv:2510.08457*.
- Xiusi Chen, Gaotang Li, Ziqi Wang, Bowen Jin, Cheng Qian, Yu Wang, Hongru Wang, Yu Zhang, Denghui Zhang, Tong Zhang, and 1 others. 2025d. [Rm-r1: Reward modeling as reasoning](#). *arXiv preprint arXiv:2505.02387*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *arXiv preprint arXiv:2110.14168*.
- Yong Deng, Guoqing Wang, Zhenzhe Ying, Xiaofeng Wu, Jinzhen Lin, Wenwen Xiong, Yuqin Dai, Shuo Yang, Zhanwei Zhang, Qiwen Wang, and 1 others. 2025. [Atom-searcher: Enhancing agentic deep research via fine-grained atomic thought reward](#). *arXiv preprint arXiv:2508.12800*.
- Guanting Dong, Hangyu Mao, Kai Ma, Licheng Bao, Yifei Chen, Zhongyuan Wang, Zhongxia Chen, Jiazhen Du, Huiyang Wang, Fuzheng Zhang, and 1 others. 2025. [Agentic reinforced policy optimization](#). *arXiv preprint arXiv:2507.19849*.
- Kaixuan Fan, Kaituo Feng, Haoming Lyu, Dongzhan Zhou, and Xiangyu Yue. 2025. [Sophiavl-r1: Reinforcing mllms reasoning with thinking reward](#). *arXiv preprint arXiv:2505.17018*.
- Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Junfei Wu, Xiaoying Zhang, Benyou Wang, and Xiangyu Yue. 2025a. [Video-r1: Reinforcing video reasoning in mllms](#). *arXiv preprint arXiv:2503.21776*.
- Kaituo Feng, Manyuan Zhang, Hongyu Li, Kaixuan Fan, Shuang Chen, Yilei Jiang, Dian Zheng, Peiwen Sun, Yiyuan Zhang, Haoze Sun, and 1 others. 2025b. [Onethinker: All-in-one reasoning model for image and video](#). *arXiv preprint arXiv:2512.03043*.
- Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. 2025c. [Group-in-group policy optimization for llm agent training](#). *arXiv preprint arXiv:2505.10978*.
- 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.
- Ziyu Hu, Zhengliang Shi, Minghang Zhu, Haitao Li, Teng Sun, Pengjie Ren, Suzan Verberne, and Zhaochun Ren. 2025. [Openreward: Learning to reward long-form agentic tasks via reinforcement learning](#). *arXiv preprint arXiv:2510.24636*.
- Yuxiang Ji, Ziyu Ma, Yong Wang, Guanhua Chen, Xiangxiang Chu, and Liaoni Wu. 2025. [Tree search for llm agent reinforcement learning](#). *arXiv preprint arXiv:2509.21240*.

- Ai Jian, Jingqing Ruan, Xing Ma, Dailin Li, Qian-Lin Zhou, Ke Zeng, and Xunliang Cai. 2025. Patarm: Bridging pairwise and pointwise signals via preference-aware task-adaptive reward modeling. *arXiv preprint arXiv:2510.24235*.
- Dongfu Jiang, Yi Lu, Zhuofeng Li, Zhiheng Lyu, Ping Nie, Haozhe Wang, Alex Su, Hui Chen, Kai Zou, Chao Du, and 1 others. 2025. Verltool: Towards holistic agentic reinforcement learning with tool use. *arXiv preprint arXiv:2509.01055*.
- 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*.
- Hongyu Li, Manyuan Zhang, Dian Zheng, Ziyu Guo, Yimeng Jia, Kaituo Feng, Hao Yu, Yexin Liu, Yan Feng, Peng Pei, and 1 others. 2025a. Editthinker: Unlocking iterative reasoning for any image editor. *arXiv preprint arXiv:2512.05965*.
- Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, and 1 others. 2025b. Websailor: Navigating super-human reasoning for web agent. *arXiv preprint arXiv:2507.02592*.
- Renhao Li, Jianhong Tu, Yang Su, Hamid Alinejad-Rokny, Derek F Wong, Junyang Lin, and Min Yang. 2025c. One model to critique them all: Rewarding agentic tool-use via efficient reasoning. *arXiv preprint arXiv:2510.26167*.
- Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025d. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*.
- Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yongkang Wu, Ji-Rong Wen, Yutao Zhu, and Zhicheng Dou. 2025e. Webthinker: Empowering large reasoning models with deep research capability. *arXiv preprint arXiv:2504.21776*.
- Zhuofeng Li, Haoxiang Zhang, Seungju Han, Sheng Liu, Jianwen Xie, Yu Zhang, Yejin Choi, James Zou, and Pan Lu. 2025f. In-the-flow agentic system optimization for effective planning and tool use. *arXiv preprint arXiv:2510.05592*.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.
- Minhua Lin, Zongyu Wu, Zhichao Xu, Hui Liu, Xianfeng Tang, Qi He, Charu Aggarwal, Xiang Zhang, and Suhang Wang. 2025. A comprehensive survey on reinforcement learning-based agentic search: Foundations, roles, optimizations, evaluations, and applications. *arXiv preprint arXiv:2510.16724*.
- Xiaoqian Liu, Ke Wang, Yuchuan Wu, Fei Huang, Yongbin Li, Junge Zhang, and Jianbin Jiao. 2025a. Agentic reinforcement learning with implicit step rewards. *arXiv preprint arXiv:2509.19199*.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025b. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*.
- Xing Han Lù, Amirhossein Kazemnejad, Nicholas Meade, Arkil Patel, Dongchan Shin, Alejandra Zambrano, Karolina Stańczak, Peter Shaw, Christopher J Pal, and Siva Reddy. 2025. Agentrewardbench: Evaluating automatic evaluations of web agent trajectories. *arXiv preprint arXiv:2504.08942*.
- Grégoire Mialon, Clémentine Fourier, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*.
- Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, and 1 others. 2025. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5687–5711.
- Salman Rahman, Sruthi Gorantla, Arpit Gupta, Swastik Roy, Nanyun Peng, and Yang Liu. 2025. Spark: Stepwise process-aware rewards for reference-free reinforcement learning. *arXiv preprint arXiv:2512.03244*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Guangming Sheng, Chi Zhang, Zilingfeng 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*.
- Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *arXiv preprint arXiv:2503.05592*.
- Sijun Tan, Michael Luo, Colin Cai, Tarun Venkat, Kyle Montgomery, Aaron Hao, Tianhao Wu, Arnab Balyan, Manan Roongta, Chenguang Wang, Li Er-ran Li, Raluca Ada Popa, and Ion Stoica. 2025. rllm: A framework for post-training language agents. <https://pretty-radio-b75.notion.site/rLL>

- M-A-Framework-for-Post-Training-Language-Agents-21b81902c146819db63cd98a54ba5f31. Notion Blog.
- Xinyu Tang, Yuliang Zhan, Zhixun Li, Wayne Xin Zhao, Zhenduo Zhang, Zujie Wen, Zhiqiang Zhang, and Jun Zhou. 2025. Rethinking sample polarity in reinforcement learning with verifiable rewards. *arXiv preprint arXiv:2512.21625*.
- 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*.
- Binghai Wang, Rui Zheng, Lu Chen, Yan Liu, Shihan Dou, Caishuang Huang, Wei Shen, Senjie Jin, Enyu Zhou, Chenyu Shi, and 1 others. 2024. Secrets of rlhf in large language models part ii: Reward modeling. *arXiv preprint arXiv:2401.06080*.
- Chaoyang Wang, Kaituo Feng, Dongyang Chen, Zhongyu Wang, Zhixun Li, Sicheng Gao, Meng Meng, Xu Zhou, Manyuan Zhang, Yuzhang Shang, and 1 others. 2025. Adatooler-v: Adaptive tool-use for images and videos. *arXiv preprint arXiv:2512.16918*.
- Chenxi Whitehouse, Tianlu Wang, Ping Yu, Xian Li, Jason Weston, Ilya Kulikov, and Swarnadeep Saha. 2025. J1: Incentivizing thinking in llm-as-a-judge via reinforcement learning. *arXiv preprint arXiv:2505.10320*.
- 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*.
- Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, and 1 others. 2025b. Webwalker: Benchmarking llms in web traversal. *arXiv preprint arXiv:2501.07572*.
- Junfei Wu, Jian Guan, Kaituo Feng, Qiang Liu, Shu Wu, Liang Wang, Wei Wu, and Tieniu Tan. 2025c. Reinforcing spatial reasoning in vision-language models with interwoven thinking and visual drawing. *arXiv preprint arXiv:2506.09965*.
- Zhiheng Xi, Chenyang Liao, Guanyu Li, Yajie Yang, Wenxiang Chen, Zhihao Zhang, Binghai Wang, Senjie Jin, Yuhao Zhou, Jian Guan, and 1 others. 2025. Agentprm: Process reward models for llm agents via step-wise promise and progress. *arXiv preprint arXiv:2511.08325*.
- Peng Xia, Kaide Zeng, Jiaqi Liu, Can Qin, Fang Wu, Yiyang Zhou, Caiming Xiong, and Huaxiu Yao. 2025. Agent0: Unleashing self-evolving agents from zero data via tool-integrated reasoning. *arXiv preprint arXiv:2511.16043*.
- Peiran Xu, Zhuohao Li, Xiaoying Xing, Guannan Zhang, Debiao Li, and Kunyu Shi. 2025. Hybrid reward normalization for process-supervised non-verifiable agentic tasks. *arXiv preprint arXiv:2509.25598*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Xiaoying Zhang, Hao Sun, Yipeng Zhang, Kaituo Feng, Chaochao Lu, Chao Yang, and Helen Meng. 2025a. Critique-grpo: Advancing llm reasoning with natural language and numerical feedback. *arXiv preprint arXiv:2506.03106*.
- Yi-Fan Zhang, Xingyu Lu, Xiao Hu, Chaoyou Fu, Bin Wen, Tianke Zhang, Changyi Liu, Kaiyu Jiang, Kaibing Chen, Kaiyu Tang, and 1 others. 2025b. R1-reward: Training multimodal reward model through stable reinforcement learning. *arXiv preprint arXiv:2505.02835*.
- Yifan Zhang and Team Math-AI. 2024. American invitational mathematics examination (aime) 2024.
- Yifan Zhang and Team Math-AI. 2025. American invitational mathematics examination (aime) 2025.
- Zheng Zhang, Ziwei Shan, Kaitao Song, Yexin Li, and Kan Ren. 2025c. Linking process to outcome: Conditional reward modeling for llm reasoning. *arXiv preprint arXiv:2509.26578*.
- Zhiwei Zhang, Hui Liu, Xiaomin Li, Zhenwei Dai, Jingying Zeng, Fali Wang, Minhua Lin, Ramraj Chandradevan, Zhen Li, Chen Luo, and 1 others. 2025d. Bradley-terry and multi-objective reward modeling are complementary. *arXiv preprint arXiv:2507.07375*.
- Zijing Zhang, Ziyang Chen, Mingxiao Li, Zhaopeng Tu, and Xiaolong Li. 2025e. R1vmr: Reinforcement learning with verifiable meta-reasoning rewards for robust long-horizon agents. *arXiv preprint arXiv:2507.22844*.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

A Dataset Details

All datasets will be publicly released to support future research.

A.1 Agent Training Data Distribution

We collect 709k question-answer pairs from publicly available datasets as our RL dataset, **Reagent-RL-709K**. The detailed distribution information is shown in Figure 1 (bottom). We randomly select 100k data from RL dataset and utilize DeepSeekV3.1 to collect problem solving trajectories with our 6 tools. The trajectories reach to the correct final answer is saved as the SFT dataset in Figure 1 (top). In total, we collect 55.6k high quality trajectories for SFT training, denoted as **Reagent-SFT-55.6K**.

A.2 Dataset Selection and Filtering

DeepMath. Each question in DeepMath is accompanied by three independently generated solutions. We remove samples for which the final answers are inconsistent across the three solutions, as such cases introduce ambiguity in supervision.

DeepScaleR. We filter out samples whose provided solutions produce answers inconsistent with the labeled ground-truth answer.

SimpleRL-Zoo. To encourage non-trivial reasoning behavior, we subsample questions by difficulty, retaining a higher proportion of medium- and hard-level questions and fewer easy ones.

MMK12. We select samples where the visual input consists of charts or tables, which can be reliably processed using OCR-based tool assistance.

PixelReasoner. We select questions that require extracting information from textual content or visual avatars within images. These samples are solvable using a combination of OCR and image description tools.

LiveVQA. We retain questions that ask about identifiable attributes such as titles or authors present in images. Only samples from the image-based subset are included.

ToolVQA. We select questions from the ImageDescription, GoogleSearch, OCR, and Calculator categories to align with the agent’s available tool set.

SimpleDeepSearcher. We convert the original tool-calling format into a Qwen-compatible format to ensure consistency with the agent’s action space.

AFM-WebAgent. We transform multi-agent interaction data into a single-agent reasoning trajectory by linearizing planning, verification, and reflection steps. These components provide useful reasoning patterns for single-agent reasoning.

LongAudio. We select audio samples with durations between 5 and 40 seconds. Audio clips shorter than 5 seconds typically lack sufficient informational content, while longer clips impose excessive computational overhead on the audio-to-text model (whisper-large-v3).

A.3 Agent-RRM Construction Details

The prompt template employed to generate training data for Agent-RRM is detailed in Figure 4. Our primary design objective is to augment the reasoning-driven analytical capabilities of Agent-RRM, ensuring it provides reliable and informative feedback at both the semantic and scalar levels. To this end, each training instance is structured into three distinct components to facilitate multi-granular reasoning supervision:

- **<think>**: A reasoning process that evaluates the logical consistency of the given trajectory. By explicitly articulating the rationale behind the assessment, this component provides transparency into how the reward model derives its judgments, thereby enhancing the interpretability of the final reward signal.
- **<critique>**: A targeted summary that identifies reasoning flaws, with a particular focus on global logic and the appropriate invocation of external tools.
- **<score>**: A holistic scalar value derived from the preceding analysis to quantify the trajectory’s overall quality. By condensing complex reasoning evaluations into a standardized numerical format, this component serves as the formal reward signal required for advantage calculation within the reinforcement learning optimization loop.

To support the development of our model, we curate two specialized datasets: **Reagent-RRM-SFT-28K**, comprising 28,000 high-quality trajectories for initial supervised fine-tuning, and **Reagent-RRM-RL-90K**, consisting of 90,000 instances designed for large-scale RL training.

Instruction Prompt for Reward Model Annotation

You are an expert agent tool use evaluator. You must strictly follow the output format below:

<think>

Provide a comprehensive analysis of the entire reasoning trajectory. Focus specifically on the agent’s reasoning quality and its tool-usage behavior across ANY type of tool.

Key points to evaluate (for all tasks and all tools):

- Whether the agent correctly decided when to call tools. Over-reliance on tools for trivial reasoning is bad; failing to call tools when necessary is also bad.
- Whether the agent misused tools (e.g., calling an irrelevant tool, giving incorrectly formatted arguments, hallucinating tool inputs or filenames, making repeated tool calls without new purpose).
- Whether the agent understood tool limitations (e.g., tool outputs may be incomplete, noisy, or partial; tools cannot access nonexistent resources).
- Whether the agent improved its reasoning over time (e.g., corrected wrong assumptions, avoided repeated mistakes, verified hypotheses when possible).
- Whether the agent avoided unverified guesses. Hypotheses without verification are harmful.
- Whether the agent avoided fabricating tool results, file names, object identifiers, or other non-existent content.

If uncertain, identify potential harmful reasoning patterns: unnecessary tool calls, missing essential tool calls, uncritical acceptance of tool output, faulty logical jumps, or incorrect assumptions about tool capabilities.

Never mention the true answer. Only evaluate the reasoning process and tool use.

</think>

<critique>

Provide a succinct, specific, and actionable summary of issues in the agent’s reasoning and tool use. This section will be shown to the agent, so it must be concise and clearly highlight:

- Incorrect, unnecessary, missing, or repeated tool calls.
- Incorrect assumptions, unverified reasoning, or blind trust in tool results.
- Any improper handling of tool limitations or constraints.
- Any hallucinated tool arguments, filenames, or resource identifiers.
- Unlogical reasoning.

Do NOT provide the correct answer or hints toward it.

</critique>

<score>

A single float between 0 and 1 representing the overall quality of the reasoning and tool use. 0 means completely incorrect or harmful reasoning; 1 means flawless reasoning with appropriate, precise, and well-justified tool use. </score>

Strict Requirements: Output exactly three blocks; focus solely on reasoning/tool-use; never reveal the correct answer.

Figure 4: The prompt used for generating structured judgments of reward model.

B Training Details

B.1 Training Codebase

We use LLaMA-Factory (Zheng et al., 2024) to implement SFT training of both Agent-RRM and Reagent. Both reward model and agent model are trained for 2 epoches.

We use rLLM (Tan et al., 2025) to implement Agentic RL training of Reagent. We use VeRL (Sheng et al., 2024) to implement RL training of Agent-RRM. The hyper-parameters we used during RL training is shown in Table 4. We conduct RL training for 300 steps.

B.2 Tools configuration

Our models are trained on 8 NVIDIA A800-80G GPUs. We detail the specific implementations of tools integrated into our agentic framework during

Table 4: Hyper-parameters for Reinforcement Learning Training.

Category	Hyper-parameter	Value
Training Config	Base Model	Qwen3-8B
	Optimizer	AdamW
	Learning Rate	5×10^{-7}
	Training Batch Size	64
	Mini-batch Size	16
Generation & Env	Max Agent Steps	13
	Temperature	0.7
	Top- p	0.95
	Rollout Samples (n)	8
	Lambda (λ)	0.3

training:

- **Search:** Powered by the Bing Search API. The agent receives the top- k results, including the URL, title, and a content snippet for each entry.
- **Browse:** Website content is retrieved via the Jina Reader and subsequently condensed using DeepSeek-Chat as a summarization model.
- **Image2text:** Visual queries and image-based reasoning are handled by GPT-4.1.
- **Audio2text:** Audio inputs are transcribed into text using the Whisper-large-v3 model.

C Evaluation Details

C.1 Evaluation Benchmarks

For GAIA, we report performance on the 103-task Text subset in Table 1 following Dong et al. (2025), while the full set results (165 tasks) are reported as GAIA (Full). For HLE, we evaluate on the 500-task subset consistent with Dong et al. (2025), as shown in Table 1. For xbench, we report results specifically on the xbench-DeepSearch set.

C.2 Tool Configurations

The configurations for evaluation tools are identical to those detailed in Appendix B.2.

C.3 Agent Evaluation Settings

For all benchmark evaluations, we employ a decoding temperature of 0.6 and a top_p of 0.95. The maximum response length is constrained to 32,768 tokens, and the agent is permitted a maximum of 30 tool-use steps per query. Consistent with Dong et al. (2025), we utilize Qwen2.5-72B-Instruct as the automated judge to perform binary scoring by comparing the agent’s predictions against the ground-truth answers.

D Case Study

We present two case studies of Reagent-C in Figure 5 and Figure 6 to demonstrate the effectiveness of our critique mechanism. The critiques generated by Agent-RRM highlight logical inconsistencies or inappropriate tool usage in the initial responses. These critiques serve as informative signals that facilitate the agent’s self-correction, leading to successful reasoning refinement and error rectification in the subsequent responses.

Question: On July 15, 2008, Phys.org published an article about a catastrophe. Find the explosive force of this catastrophe according to Encyclopedia Britannica, then find the name of the US nuclear test that had the same yield. Your answer should only be the last word of the name of the test.

Ground Truth: Bravo

Source: GAIA

First Response:

<think>

Okay, let's tackle this step by step. The user wants the explosive force from a Phys.org article on July 15, 2008, then find a US nuclear test with the same yield... I'll search for the article first.

</think>

<tool_call>

{"name": "search", "arguments": {"queries": ["Phys.org July 15 2008 catastrophe"]}}

</tool_call>

<tool_response>

--- search result for [Phys.org July 15 2008 catastrophe] ---

[Tunguska Catastrophe: Evidence Of Acid Rain...] The Tunguska catastrophe in 1908 evidently...

[Planetary science: Tunguska at 100 - Nature] ...it was a 15-megatonne event...

[Was there an upward atmospheric discharge...] ...the Tunguska explosion ...explosion energy of 3.5–5.0 megatons...

</tool_response>

<think>

...The search results mention both a 15-megatonne event and 3.5–5.0 megatons. The Britannica entry might cite a specific value. Let's assume the yield is around 5 megatons based on the range provided. Now I need to find a US nuclear test with that yield...

</think>

<tool_call>

{"name": "search", "arguments": {"queries": ["US nuclear test 5 megaton yield"]}}

</tool_call>

<tool_response>

--- search result for [US nuclear test 5 megaton yield] ---

[The BRAVO Test - Atomic Archive] The predicted was 5 megatons, but "BRAVO" was 14.8 megatons...

[Cannikin Nuclear Test Site] ...consisted of a 5 megaton-yield thermonuclear bomb...

</tool_response>

<think>

...The search results show the Tunguska event's yield is 3-5 megatons. The Britannica entry mentions 185 Hiroshima bombs, roughly equivalent to 5 megatons. So the "Cannikin" test had a 5-megaton yield. Therefore, the last word is **Cannikin**.

</think>

Critique:

- **Missing essential tool call:** The agent never used browse tool to open the Phys.org article or the Britannica page, so it never confirmed the catastrophe's name or the exact yield.

- **Unverified assumptions:** The agent inferred the catastrophe was the Tunguska event and the yield was 5 megatons without extracting information directly from the Britannica entry.

- **Ambiguous test selection:** The agent listed multiple US tests (Cannikin, Castle Bravo, Buster-Jangle) but chose "Cannikin" without checking which the Britannica yield matches.

Second Response:

<think>

...Following the critique, I should use the browse tool to get the exact yield from Britannica...

</think>

<tool_call>

{"name": "browse", "arguments": {"url": "https://www.britannica.com/event/Tunguska-event", "query": "What is the explosive force of the Tunguska event in megatons?"}}

</tool_call>

<tool_response>

Based on the source content, the explosive force of the Tunguska event is estimated to have been equivalent to as much as 15 megatons of TNT...

</tool_response>

<think>

Now I have the exact yield from Britannica: 15 megatons...The last word of "Castle Bravo" is **Bravo**.

</think>

Figure 5: Case 1: a search question from GAIA.

Question: A team of 4 painters worked on a mansion for $\frac{3}{8}$ ths of a day every day for 3 weeks. How many hours of work did each painter put in?

Ground Truth: 189

Source: GSM8K

First Response:

<think>

Okay, let's try to figure out how many hours each painter worked...First, I need to break down the time they worked. Let me start by calculating..., which is 21 days. So 3 weeks = 21 days. Next, each day they worked $\frac{3}{8}$ of a day...

</think>

<tool_call>

```
{"name": "python", "arguments": {"code": "total_weeks = 3\nwork_days_per_week = 7\ndays_worked = total_weeks * work_days_per_week\nfraction_of_day = 3/8\ntotal_days = days_worked * fraction_of_day\nhours_per_day = 24\ntotal_hours = total_days * hours_per_day\npainters = 4\nhours_perPainter = total_hours / painters\nhours_perPainter"}}
```

</tool_call>

<tool_response>

Code executed successfully (no output)

</tool_response>

<think>

Okay, let me check... total days worked by the team is $21 * \frac{3}{8} = \frac{63}{8}$ days. Convert that to hours by multiplying by 24: $\frac{63}{8} * 24 = 63 * 3 = 189$ hours. Divide by 4 painters: $189 / 4 = 47.25$ hours per painter. The Python code calculate this...

</think>

<tool_call>

```
{"name": "python", "arguments": {"code": "total_weeks = 3\nwork_days_per_week = 7\ndays_worked = total_weeks * work_days_per_week\nfraction_of_day = 3/8\ntotal_days = days_worked * fraction_of_day\nhours_per_day = 24\ntotal_hours = total_days * hours_per_day\ntotal_hours_perPainter = total_hours / 4\nprint(total_hours_perPainter)"}}
```

</tool_call>

<tool_response>

47.25

</tool_response>

<think>

Okay, let me ...So total days worked by the team is $21 * \frac{3}{8} = \frac{63}{8}$ days... $\frac{63}{8} * 24 = 189$ hours. Divided by 4 painters, that's $189 / 4 = 47.25$ hours per painter. The Python code correctly calculates this, and the output is 47.25. So the answer should be **47.25 hours**.

</think>

Critique:

- The first Python call did not include a `print` statement, so the tool returned no output. This is an unnecessary tool call because the assistant already had the answer, and it wasted a call.

- The assistant should have either omitted the first tool call or included a `print` to capture the result.

Second Response:

<think>

Let me carefully work through this problem step by step. Given:

- Team of 4 painters ...

Step 1: Convert $\frac{3}{8}$ of a day to hours

1 day = 24 hours $\frac{3}{8}$ of a day = $(\frac{3}{8}) \times 24 = 9$ hours per day

Step 2: Calculate total days: 3 weeks = $3 \times 7 = 21$ days

Step 3: Calculate hours per painter

Since EACH painter worked 9 hours per day for 21 days: Hours per painter = $9 \times 21 = 189$ hours

Wait, should I divide by 4? Let me think about this carefully.

The question asks for how many hours EACH painter worked. If each painter worked alongside the others for 9 hours/day for 21 days, then EACH painter individually put in 189 hours of work.

Total man-hours = $4 \times 189 = 756$ hours, but the question asks for EACH painter, which is **189 hours**.

</think>

Figure 6: Case 2: a math question from GSM8K.