
Agent-R1: Training Powerful LLM Agents with End-to-End Reinforcement Learning

Mingyue Cheng, Jie Ouyang, Shuo Yu, Ruiran Yan, Yucong Luo,
Zirui Liu, Daoyu Wang, Qi Liu, Enhong Chen
State Key Laboratory of Cognitive Intelligence,
University of Science and Technology of China, Hefei, China

Abstract

Large Language Models (LLMs) are increasingly being explored for building Agents capable of active environmental interaction (e.g., via tool use) to solve complex problems. Reinforcement Learning (RL) is considered a key technology with significant potential for training such Agents; however, the effective application of RL to LLM Agents is still in its nascent stages and faces considerable challenges. Currently, this emerging field lacks in-depth exploration into RL approaches specifically tailored for the LLM Agent context, alongside a scarcity of flexible and easily extensible training frameworks designed for this purpose. To help advance this area, this paper first revisits and clarifies Reinforcement Learning methodologies for LLM Agents by systematically extending the Markov Decision Process (MDP) framework to comprehensively define the key components of an LLM Agent. Secondly, we introduce Agent-R1, a modular, flexible, and user-friendly training framework for RL-based LLM Agents, designed for straightforward adaptation across diverse task scenarios and interactive environments. We conducted experiments on Multihop QA benchmark tasks, providing initial validation for the effectiveness of our proposed methods and framework.

🔗 <https://github.com/0russwest0/Agent-R1>.¹

1 Introduction

In recent years, large language models (LLMs) have demonstrated remarkable capabilities in natural language understanding and generation [2, 4], and are increasingly being applied to more complex intelligent tasks[24]. When LLMs are assigned an “agent” role, they are expected not only to perform cognitive tasks such as reasoning and decision-making, but also to act autonomously, learn continuously, and adapt to changes within interactive environments [35, 25, 16]. Unlike traditional static reasoning tasks[30], LLMs functioning as agents must maintain memory across multiple dialogue rounds[37], exhibit sequential decision-making capabilities, and respond effectively to environmental feedback—bringing them closer to real-world autonomous intelligent systems [27, 13]. This direction opens new possibilities for building general artificial intelligence with self-evolving and problem-solving capabilities[31, 9].

While Reinforcement Learning (RL)[5, 26, 6] has demonstrated notable success in enhancing LLM capabilities for relatively well-defined tasks such as mathematical problem solving and code generation [15, 19, 1], its application to developing LLMs as autonomous, interactive agents is comparatively nascent. Agent settings inherently require models to make sequential decisions, maintain memory across turns, and adapt to stochastic environmental feedback [35, 27], presenting unique challenges distinct from more static tasks[16, 30]. This leads to specific difficulties when applying RL; particularly in multi-turn interaction scenarios, agent training can encounter instability, complex reward

¹This paper serves as the technical report of the Agent-R1 project.

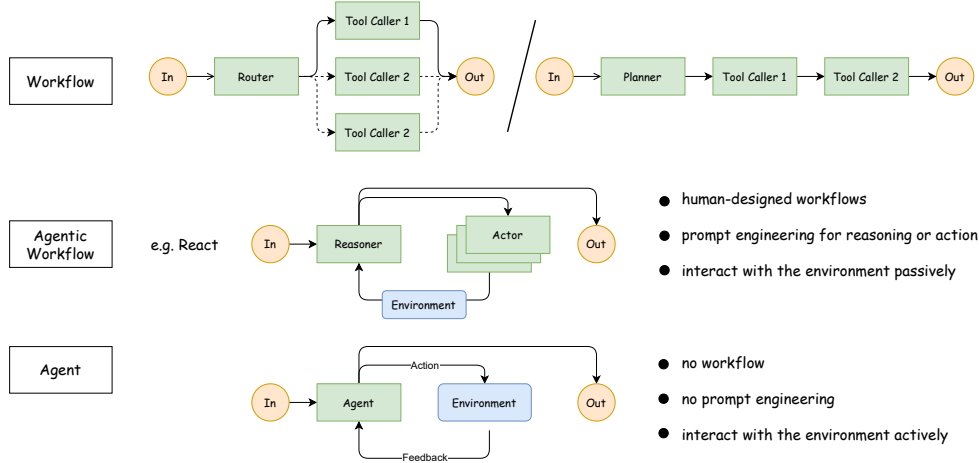


Figure 1: Comparison of workflows, agentic workflows, and autonomous agents. Workflows rely on human-designed routing or planning, while agentic workflows (e.g., ReAct) introduce iterative reasoning-acting loops. Fully autonomous agents remove predefined workflows and interact with the environment proactively through an end-to-end action-feedback cycle.

signal design, and limited generalization [3, 10, 20]. Thus, there remains a need for a more detailed exploration of how RL methodologies can be systematically applied and adapted to LLM agents, alongside a need for flexible and scalable training frameworks [14, 18].

To systematically address these aspects, this paper presents contributions from both a conceptual and practical standpoint. Conceptually, we focus on clarifying the application of Reinforcement Learning to LLM Agents. We achieve this by extending the standard Markov Decision Process (MDP) framework[22], providing a detailed account of how its core components—State Space, Action Space, State Transition Probability, and Reward Function—can be adapted to comprehensively model the multi-turn, interactive nature of LLM Agents [7, 36]. Building upon this adapted MDP formulation, we further elaborate on mechanisms for optimizing agent policy from multi-turn trajectories[17], emphasizing the importance of distinguishing agent-generated actions from environmental feedback and incorporating intermediate (process) rewards [11, 27] to guide learning effectively. In addition, to facilitate the practical application of these concepts, we develop **Agent-R1**, a flexible and user-friendly training platform for RL-based LLM Agents. Thanks to its modular architecture, Agent-R1 supports rapid integration of various environment interfaces and task scenarios, and can dynamically adapt to different computational resource requirements, enabling easy scalability to complex and diverse applications [32, 12].

Through systematic experiments on the challenging task of multi-hop question answering (Multi-hop QA), we validate the effectiveness of our approach and framework[16, 33]. This task focuses on complex reasoning involving cross-document logical chaining and information retrieval, imposing high demands on the agent’s multi-step decision-making ability, adaptability to environmental feedback, and knowledge construction processes. Experimental results demonstrate that our methods and framework can enhance the model’s performance in such dynamic interactive environments[29].

2 From Large Language Models to Agents: An MDP Perspective

The sequential decision-making processes inherent in Large Language Model (LLM) applications can be effectively formulated within the Markov Decision Process (MDP) framework, whether for simple text generation or for the complex interactions of LLM Agents. However, evolving the MDP formulation from its application to LLMs in the context of **static, single-turn** text generation tasks (such as in mathematics or code generation), to one suitable for LLM Agents that engage in inherently **dynamic, multi-turn**, and richly interactive environmental dialogues, necessitates substantial extensions. This section delineates these critical differences by contrasting the core MDP components for a **Static LLM** versus a **LLM Agent**.

Table 1: Comparison of Markov Decision Process (MDP) components for a Static LLM versus an LLM Agent, highlighting the necessary extensions for interactive, multi-turn scenarios.

MDP Component	Static LLM	LLM Agent
State Space (\mathcal{S})	The state s_t primarily encapsulates the current textual context, including the initial prompt and the sequence of tokens generated so far. Its focus is on predicting the next coherent token. $s_t = (w_p, w_1, \dots, w_t)$.	The state s_t is significantly more comprehensive, retaining the history of multi-turn interactions and environmental feedback. This enables decisions informed by a full dialogue history. Each \mathcal{T}_i represents a full turn of agent action and environment feedback. $s_t = (w_p, \mathcal{T}_1, \dots, \mathcal{T}_k, \mathcal{T}_{k+1}^{\text{partial}})$.
Action Space (\mathcal{A})	The action a_t corresponds to selecting the next token w_{t+1} from the LLM’s vocabulary V .	The fundamental action is also token generation from V . However, specific token sequences can be interpreted as commands to invoke external tools, enabling active environmental intervention.
State Transition (\mathcal{P})	State transitions are deterministic. The next state is uniquely determined by appending the selected token to the current sequence. $P(s_{t+1} s_t, a_t) = 1$ if $s_{t+1} = s_t \oplus a_t$.	The state transition mechanism incorporates environmental interaction, which can be stochastic. It distinguishes between deterministic generative transitions (P_G) and potentially stochastic environmental transitions (P_E) triggered by tool use.
Reward Function (\mathcal{R})	Rewards are typically sparse and provided at the end of a complete generation, evaluating the overall quality of the final output.	The reward structure is richer and more dense. In addition to a final outcome reward (r_f), agents can receive intermediate process rewards (r_p) for successfully executing steps like tool invocation, providing more frequent feedback.

2.1 State Space (\mathcal{S})

Static LLM: In single-turn text generation, the state s_t primarily encapsulates the current textual context. This includes the initial prompt w_p and the sequence of tokens generated thus far w_1, w_2, \dots, w_t :

$$s_t = (w_p, w_1, w_2, \dots, w_t). \quad (1)$$

The state space is focused on capturing the information necessary to predict the next token in a coherent sequence.

LLM Agent: For an Agent engaging in multi-turn interactions, the state s_t must be significantly more comprehensive. It needs to retain not only the textual context but also the history of interactions and environmental feedback. The state is thus extended to:

$$s_t = (w_p, \mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_k, \mathcal{T}_{k+1}^{\text{partial}}). \quad (2)$$

Here, each \mathcal{T}_i represents a complete interaction turn, comprising the Agent’s generated tokens $(w_{i1}, \dots, w_{iT_i})$ and the subsequent environmental feedback w_{ei} , i.e., $\mathcal{T}_i = (w_{i1}, \dots, w_{iT_i}, w_{ei})$. The term $\mathcal{T}_{k+1}^{\text{partial}}$ denotes the partially generated sequence in the current, ongoing turn. This enriched state representation enables the Agent to make decisions informed by a comprehensive history of dialogue and environmental outcomes (e.g., results from tool utilization).

Table 2: Core differences in the MDP formulation for Static LLMs vs. LLM Agents.

Component	Static LLM	LLM Agent
State (S)	Captures only the current text sequence.	Captures the full history of multi-turn interactions and environmental feedback.
Action (A)	Generating the next token.	Generating tokens that can also function as commands to invoke external tools.
State Transition (P)	Deterministic: appending a token determines the next state.	Stochastic: the next state depends on non-deterministic feedback from environment.
Reward (R)	Receives a single, sparse reward at the end of the generation.	Receives dense process rewards for intermediate steps in addition to a final reward.

2.2 Action Space (A)

Static LLM: The action a_t corresponds to the selection of the next token w_{t+1} from the LLM’s vocabulary \mathcal{V} . The action space $A(s_t)$ is typically the set \mathcal{V} .

LLM Agent: Similarly, an Agent’s action a_t is the selection of the next token from \mathcal{V} . However, the *implication* of an action sequence can be broader. Specific sequences of tokens generated by the Agent might be interpreted as commands to invoke external tools or APIs. Thus, while the fundamental action remains token generation, its functional outcome can extend beyond mere text production to active environmental intervention.

2.3 State Transition Probability (P)

Static LLM: State transitions in LLM text generation are **deterministic**. Given the current state s_t and an action a_t (selecting token w_{t+1}), the next state s_{t+1} is uniquely determined by appending w_{t+1} to the current sequence:

$$P(s_{t+1}|s_t, a_t) = \begin{cases} 1, & \text{if } s_{t+1} = s_t \oplus a_t \\ 0, & \text{otherwise} \end{cases}, \quad (3)$$

where \oplus denotes sequence concatenation.

LLM Agent: The state transition mechanism for Agents introduces a crucial distinction by incorporating **environmental interaction**, which can be stochastic. Transitions can be categorized based on whether an action instigates such interaction:

$$P(s_{t+1}|s_t, a_t) = \begin{cases} P_E(s_{t+1}|s_t, a_t), & \text{if } a_t \text{ triggers tool/environment interaction} \\ P_G(s_{t+1}|s_t, a_t), & \text{otherwise (standard token generation)} \end{cases}, \quad (4)$$

While P_G (generative transition) mirrors the deterministic nature of Static LLM token generation ($P_G(s_{t+1}|s_t, a_t) = 1$ if $s_{t+1} = s_t \oplus a_t$, and 0 otherwise), P_E (environmental transition) reflects the uncertainty inherent in tool execution and environmental responses. The next state s_{t+1} then depends not only on the Agent’s action but also on the outcome from the external environment (e.g., an API response, result of a computation), which forms part of the environmental feedback w_{ei} .

2.4 Reward Function (R)

Static LLM: Rewards are typically **sparse** and provided at the end of a complete generation sequence, i.e., upon reaching a terminal state s_T . This is often a result-based reward, $R(s_T)$, evaluating the overall quality (e.g., coherence, relevance) of the generated text.

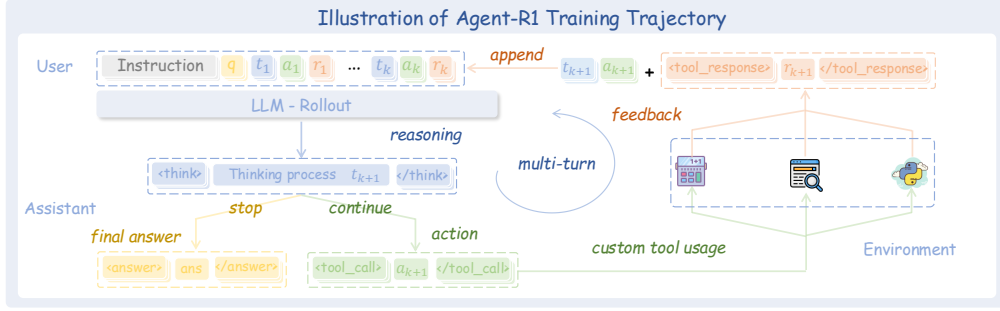


Figure 2: Illustration of the Agent-R1 training trajectory. The agent performs multi-turn reasoning and tool-based actions during rollout, receives environment feedback, and appends tool responses to form the next state. This trajectory—containing thinking steps, actions, and feedback—serves as the basis for reinforcement learning updates in Agent-R1.

LLM Agent: The reward structure for Agents is often richer and more **dense**, accommodating the multi-turn nature of their tasks. The reward $R(s_t, a_t, s_{t+1})$ can be defined as:

$$R(s_t, a_t, s_{t+1}) = \begin{cases} r_f(s_{t+1}), & \text{if } s_{t+1} \text{ is a terminal state} \\ r_p(s_t, a_t, s_{t+1}), & \text{if } a_t \text{ triggers a significant intermediate event} \\ 0, & \text{otherwise (e.g., during routine token generation)} \end{cases} \quad (5)$$

Here, $r_f(s_{t+1})$ is the final **outcome reward** for task completion. Crucially, Agents can also receive **process rewards**, $r_p(s_t, a_t, s_{t+1})$, for successfully executing intermediate steps, such as effective tool invocation or making tangible progress towards a goal. These intermediate signals provide more frequent feedback, guiding the learning process more effectively.

In summary, adapting the MDP framework from Static LLMs to LLM Agents involves significant enhancements across the core components discussed. The state space expands to incorporate interaction history and environmental feedback; actions, while fundamentally token generation, can trigger external effects; state transitions integrate environmental stochasticity; and the reward system becomes more granular with the inclusion of process rewards. These extensions are crucial for enabling reinforcement learning algorithms to train sophisticated Agents capable of complex, multi-step reasoning and interaction within dynamic environments.

3 Agent-R1 Framework

To better accommodate the reinforcement learning requirements of LLM Agents, we introduce **Agent-R1**, a flexible and highly extensible agent reinforcement learning training framework as shown in Figure 2. Drawing from existing efficient reinforcement learning infrastructures, we extend traditional single-turn reinforcement learning training frameworks to fully adapt to the multi-turn interactive characteristics of agents, enabling seamless integration with diverse task environments and scalable training across increasingly complex agent settings. Figure 3 and Figure 4 illustrate the workflow comparison between traditional single-turn reinforcement learning training frameworks and Agent-R1’s multi-turn reinforcement learning training, where Figure 3 shows the generation stage and Figure 4 shows the learning stage. The most significant difference between single-turn and multi-turn reinforcement learning lies in the rollout phase: single-turn rollout processes only require the Actor Model to generate responses once, while multi-turn rollout involves multiple complex interactions. To achieve flexible and easily extensible multi-turn rollout, we have carefully designed two core modules: **Tool** and **ToolEnv**.

3.1 Tool and ToolEnv: Core Modules for Interactive Rollout

The interactive rollout process, central to training LLM Agents, relies heavily on two key components: **Tool** and **ToolEnv**. A clear division of responsibilities between these modules is fundamental to the **Agent-R1** design philosophy.

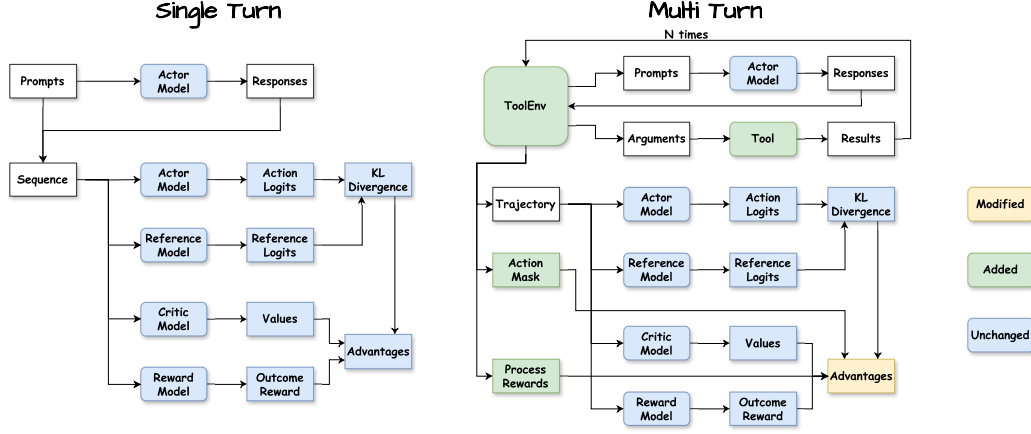


Figure 3: Flow diagram of Single-Turn RL and Multi-Turn RL(Agent-R1) in generation stage.

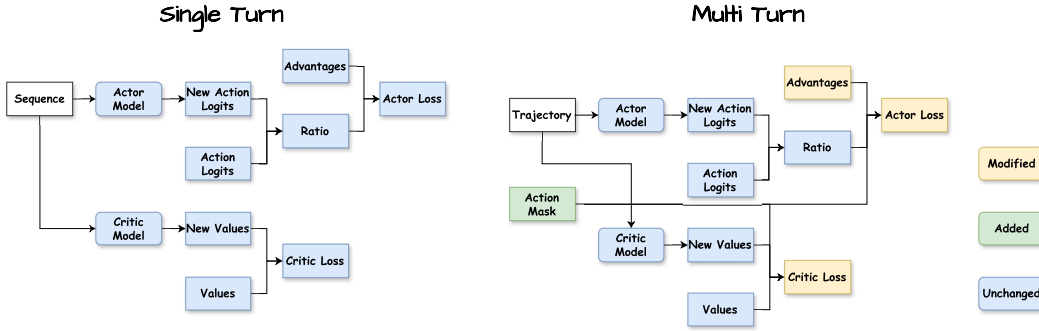


Figure 4: Flow diagram of Single-Turn RL and Multi-Turn RL(Agent-R1) in learning stage.

The **Tool** is conceived as an executor of specific, atomic actions. Its primary role is to encapsulate a distinct capability—such as calling an external API, executing a piece of code, or accessing a database. When invoked, a **Tool** performs its action and returns the direct, raw outcome of that action. It essentially reports “what happened” factually.

Conversely, the **ToolEnv** acts as the orchestrator and interpreter within the reinforcement learning (RL) environment. It takes the raw output from a **Tool** and determines how that output affects the Agent’s perceived state and the overall task progression. **ToolEnv** is responsible for managing the state transitions within the RL loop, calculating appropriate reward signals based on these transitions and tool outcomes, and packaging the new state information for the agent. It dictates “what this outcome means for the agent and the task.”

3.1.1 Tool Design

Tools serve as the critical interface connecting agents to external environments or functionalities. In the Agent-R1 framework, we utilize Tools as the unified interface for agent-environment interaction, where all external functionalities are encapsulated into standardized, directly callable “tools” by the agent. Drawing inspiration from OpenAI’s Function Calling paradigm, the Agent-R1 framework provides high-level abstraction and standardization of Tools through the BaseTool abstract base class. Its design focuses on two core modules:

Core Execution Logic: As the most critical abstract method within the BaseTool class, all concrete tool subclasses must implement the execution method. This method encapsulates the core logic of the tool, defining how it processes input parameters, performs its specific operation (such as interacting with an external API, executing code, or accessing a database), and returns a structured result.

Tool Metadata Specification: To ensure the standardization and parsability of tool invocation, the following metadata attributes are defined:

- **Identification and Description:** The name attribute (a unique string identifier) and the description attribute (providing detailed information about the tool’s functionality, use cases, and expected effects) work synergistically. Agents understand these to identify and select appropriate tools based on the current context.
- **Parameter Structure Definition:** The parameters attribute follows JSON Schema specifications to define the input parameter structure required for tool invocation. This includes parameter names, data types, detailed descriptions, and whether they are required. The standardization of parameters ensures agents can generate tool invocation parameters that conform to expected formats.

This design, centered on the execute method’s pivotal role in performing actions and the clear metadata specifications that enable agent understanding, allows LLM agents to effectively interact with external environments via structured interfaces. The outcomes of tool execution are then processed by the ToolEnv module, which is responsible for managing the corresponding environmental state transitions. This interplay is fundamental for agents to solve complex problems in multi-turn interactions and forms a cohesive link with ToolEnv’s state management design.

3.1.2 ToolEnv Design

The **ToolEnv** module acts as the dynamic environment within the Agent-R1 reinforcement learning framework. It is responsible for managing the agent’s interaction with the world, particularly when tools are involved. This module implements the two core functionalities required of an RL environment: state transition and reward calculation, especially in the context of multi-turn interactions and non-deterministic outcomes from tool use. The design is formalized through the BaseToolEnv abstract base class.

Core State Transition and Reward Logic: The **most critical** abstract method is step. This method is the primary engine for environmental interaction: it receives the agent’s raw output (e.g., generated text potentially containing tool calls) and processes this output to identify and orchestrate tool invocations in coordination with the Tool module. Based on the agent’s actions and any feedback from tool execution, step then updates the environment’s internal state. It also calculates an appropriate reward signal reflecting the outcome of the action and the new state. Finally, it returns the new state, the reward, and other relevant information (such as success status and activity flags) to the agent. This method encapsulates the logic for both standard generative state transitions and the more complex, potentially stochastic transitions resulting from tool interactions.

Supporting Mechanisms for Interaction Management: To facilitate the step method’s comprehensive role and manage the nuances of tool-based interactions, several key auxiliary methods are defined within BaseToolEnv. The process_responses_ids method provides customizable logic for identifying tool call triggers within the raw token ID sequences generated by the LLM, determining the precise point of invocation. Subsequently, extract_tool_calls is responsible for parsing these raw LLM responses to identify and structure intended tool invocation requests, including the tool’s name and its parameters. Following a tool’s execution, format_tool_response converts the raw results (obtained from Tool.execute) into a string format suitable for presentation to the LLM as part of the new environmental state. Complementing these, the stop method implements the logic for determining trajectory termination conditions, assessing whether the current interaction should end based on LLM output, task completion, error states, or predefined limits.

This design, centered on the step method’s pivotal role in driving environmental dynamics and supported by clear mechanisms for managing tool calls and trajectory lifecycle, enables the Agent-R1 framework to effectively simulate complex interactive scenarios. It carefully distinguishes between deterministic text generation and the non-deterministic, environment-altering state changes introduced by tool use, which are crucial for agent learning.

3.2 Optimizing Agent Policy from Multi-Turn Trajectories

After the rollout phase, we have obtained comprehensive multi-turn interaction trajectories. Each trajectory contains the sequence of states, the agent’s actions (generated text portions), and the reward signals. As noted, the environment provides reward signals for each interaction turn, which are termed **process rewards** (r_p), in addition to a potential final **outcome reward** (r_f). To clearly delineate the tokens generated by the LLM agent (its actions) from the environmental feedback or the initial prompt within the trajectory, we introduce an **Action Mask**. This mask identifies precisely which parts of the sequence correspond to the agent’s learnable actions.

Reinforcement learning optimizes the policy model’s actions to maximize expected cumulative rewards. The Agent-R1 framework leverages the detailed information from these multi-turn trajectories, including the action masks and process rewards, to achieve this during its learning stage (as illustrated in Figure 4). Key aspects of how this information is utilized include:

Refined and Aligned Advantage Calculation As seen in the generation stage (Figure 3), the “Advantages” are no longer solely based on the final outcome reward and value estimates from the Critic Model. The “Process Rewards” gathered from ToolEnv during the rollout are explicitly incorporated. This means the advantage \hat{A}_t at each relevant step t within a trajectory reflects not only future discounted rewards (derived from outcome rewards and value function estimates) but also the immediate success of intermediate steps, such as effective tool invocation (captured by process rewards). The “Advantages” block in Figure 3 shows it taking input from “Values” (Critic Model), “Outcome Reward” (Reward Model), and “Process Rewards”. Crucially, the calculation of these advantages \hat{A}_t (e.g., using Generalized Advantage Estimation - GAE) is performed such that they are **aligned with the agent’s actions** as identified by the Action Mask. While rewards accrue based on state transitions and value functions estimate state goodness, the final advantages used for updating the policy are pertinent to the specific timesteps where the agent generated an action. This ensures that credit (positive or negative advantage) is assigned to the actual decisions made by the agent, rather than to parts of the sequence it did not control, such as prompt tokens or fixed environmental responses. These action-aligned advantages are then passed to the learning stage (Figure 4).

Masked Policy Optimization (Actor Loss) In the learning stage (Figure 4), the **Actor Model** (policy) is updated to increase the probability of actions that lead to higher advantages. The “Trajectory” data is fed into the Actor Model to produce “New Action Logits”. The **Action Mask** plays a crucial role here. When calculating the **Actor Loss** (often a policy gradient loss like PPO’s clipped surrogate objective), the mask ensures that the loss is computed *only* over the tokens generated by the agent. The “Ratio” between the new policy’s action probabilities and the old policy’s action probabilities (derived from “New Action Logits” and “Action Logits” from the generation phase, respectively) is modulated by these aligned advantages, and this calculation is guided by the **Action Mask**.

Value Function Update (Critic Loss) The **Critic Model** is trained to more accurately estimate the expected cumulative reward (value) from different states. Using the “Trajectory” data, it generates “New Values”. The **Critic Loss** is typically a mean squared error between these “New Values” and the observed returns (which include both process and outcome rewards) from the trajectory, or the target values derived from these returns and existing value estimates (e.g., in TD learning). This helps the critic provide better baseline estimates for advantage calculation in subsequent iterations.

By ensuring that the “Advantages” are calculated in alignment with the agent’s actual actions (and subsequently used with the “Action Mask” during policy optimization), Agent-R1 provides a more precise and effective learning signal. This detailed feedback mechanism allows the Actor and Critic models to learn more efficiently from complex, extended dialogues and tool-use scenarios, driving the agent towards mastering sophisticated tasks.

4 Empirical Study

We empirically evaluate Agent-R1’s efficacy and design contributions in a challenging multi-hop question answering scenario where LLMs use external search. This study first validates the framework’s effectiveness in training LLM agents with various Reinforcement Learning (RL) algorithms for multi-turn interactive tasks. Second, an ablation analysis investigates the impact of key policy

optimization refinements: the action mask for loss computation (“loss mask”) and for aligning advantages (“advantage mask”). The overarching goal is to assess the LLM’s learned ability for tool invocation and information retrieval, highlighting Agent-R1’s utility.

4.1 Experimental Setup

Tasks and Datasets Our study uses Multi-hop Question Answering (MultihopQA) datasets. The training set comprises 51,200 samples, equally and randomly drawn from the HotpotQA[34] and 2WikiMultihopQA[8] training splits. We evaluate on the full development sets of HotpotQA and 2WikiMultihopQA (in-domain) and Musique[28] (out-of-domain), which require multi-step retrieval and reasoning.

Models and Tools Experiments use *Qwen2.5-3B-Instruct* [23] within the NousToolEnv using its native function calling. The agent employs a single wikisearch tool querying a KILT Wikipedia corpus (36M passages [21], embeddings by bge-large-en-v1.5), returning top 5 documents.

RL Algorithms and Baselines We evaluate PPO, GRPO, REINFORCE++, REINFORCE++Baseline, and RLOO to assess Agent-R1’s adaptability. These are compared against two baselines: Naive RAG (single-pass retrieval) and Base Tool Call (native function calling with wikisearch tool).

Reward Formulation A sparse final outcome reward r_f is used, defined as:

$$r_f = \begin{cases} r_{\text{answer}}, & \text{if } r_{\text{format}} = 1 \\ r_{\text{format}} - 1, & \text{if } r_{\text{format}} < 1 \end{cases}, \quad (6)$$

Here, $r_{\text{answer}} = \text{EM}(a_{\text{pred}}, a_{\text{gold}})$ is the Exact Match score. The formatting score, $r_{\text{format}} = (r_{\text{format}_a} + r_{\text{format}_t})/2$, averages binary indicators for correct final answer presentation (r_{format_a}) and valid tool call syntax (r_{format_t}). This structure strictly rewards perfectly formatted, correct answers and penalizes any formatting errors.

4.2 Main Results

The primary results of our framework validation, evaluating the performance of various Reinforcement Learning (RL) algorithms supported by Agent-R1 against baseline methods, are presented in Table 3. The experiments were conducted on three multi-hop question answering datasets. HotpotQA and 2WikiMultihopQA (2Wiki) served as in-domain datasets, while Musique was used as an out-of-domain dataset. The reported scores are Exact Match (EM) values, representing the primary metric for task success.

Table 3: Exact Match (EM) Performance Comparison of RL Algorithms and Baselines on MultihopQA Datasets. [†]Denotes in-domain datasets; ^{*}denotes out-of-domain dataset. Among RL algorithms, best performance per column is in **bold**, second best is underlined.

Method	HotpotQA [†]	2Wiki [†]	Musique [*]	Average
Base Tool Call	0.1372	0.0891	0.0277	0.0847
Naive RAG	0.1916	0.1792	0.0277	0.1328
PPO	<u>0.4136</u>	0.5468	0.1552	<u>0.3719</u>
GRPO	0.4405	0.5741	<u>0.1485</u>	0.3877
REINFORCE++	0.3768	0.4796	0.1336	0.3300
REINFORCE++Baseline	0.3966	0.5406	<u>0.1485</u>	0.3619
RLOO	0.4089	<u>0.5641</u>	0.1419	0.3716

Table 3 clearly demonstrates that all RL-trained agents substantially outperform both the Base Tool Call (0.0847) and Naive RAG (0.1328) baselines. For instance, the weakest RL agent (REINFORCE++, average EM 0.3300) still surpassed RAG by a factor of approximately 2.5. This significant margin highlights the crucial role of RL in training proficient LLM agents capable of

complex multi-turn decision-making and effective tool use, moving beyond simpler heuristic or single-pass methods.

Among the RL methods, GRPO (average EM 0.3877) exhibited the best overall performance, closely followed by PPO (0.3719) and RLOO (0.3716). PPO notably excelled on the challenging out-of-domain Musique dataset. While REINFORCE++ (0.3300) was the weakest RL performer, incorporating a baseline in REINFORCE++Baseline (0.3619) provided a clear benefit, though it did not match the top-tier algorithms. These results robustly validate Agent-R1’s efficacy in training powerful LLM agents via end-to-end RL, showing consistent, substantial gains over baselines across diverse datasets and RL algorithms. This affirms the value of our framework for optimizing agent policies in interactive settings.

4.3 Ablation Study on Policy Optimization Refinements

To investigate the significance of specific policy optimization refinements within the Agent-R1 framework—namely, the use of an action mask for loss computation (termed “loss mask”) and for the alignment of advantages (termed “advantage mask”)—we conducted an ablation study. This study utilized the PPO and GRPO algorithms. The results, in terms of Exact Match (EM) scores, are presented in Table 4, where each ablation step is shown relative to the configuration in the preceding row for that algorithm.

Table 4: Ablation Study on Policy Optimization Components. Scores are Exact Match (EM). \hookrightarrow indicates a configuration derived from the *immediately preceding row* within its algorithm group by disabling the specified component. Baseline configurations (first row of each group) use \checkmark to denote enabled components. The ‘advantage mask’ is not separately ablated for GRPO in this setup.

Configuration	HotpotQA	2Wiki	Musique	Average
PPO (loss mask \checkmark , adv. mask \checkmark)	0.4136	0.5468	0.1552	0.3719
\hookrightarrow <i>Advantage mask disabled</i>	0.3630	0.4641	0.1138	0.3136
\hookrightarrow <i>Loss mask disabled</i>	0.3429	0.4631	0.1005	0.3022
GRPO (loss mask \checkmark)	0.4405	0.5741	0.1485	0.3877
\hookrightarrow <i>Loss mask disabled</i>	0.4260	0.5485	0.1422	0.3722

The ablation experiments (Table 4) underscore the critical roles of both the loss mask and the advantage mask. Disabling the **loss mask** consistently degrades performance for both PPO (e.g., average EM dropping from 0.3136 to 0.3022 when removed from a PPO variant with the advantage mask already disabled) and GRPO (average EM dropping from 0.3877 to 0.3722). This highlights its necessity in focusing gradients on agent-generated tokens. Similarly, for PPO, disabling the **advantage mask** (while the loss mask is active, comparing the first two PPO rows) causes a substantial performance drop from an average EM of 0.3719 to 0.3136, affirming the importance of accurate credit assignment. These findings validate that these masking strategies are crucial design choices within Agent-R1 for effective policy optimization in interactive LLM agents.

5 Conclusion

This work clarifies how Reinforcement Learning can be effectively applied to LLM Agents by extending the classical MDP framework to capture multi-turn interaction, environmental feedback, and process rewards. Based on these insights, we introduced **Agent-R1**, a modular and extensible framework that supports multi-turn rollouts, precise credit assignment, and flexible integration of tools and environments. Experiments on multi-hop question answering demonstrate that Agent-R1 enables LLM agents to achieve substantial improvements over baseline methods, and ablation results confirm the importance of its key policy optimization components. We hope Agent-R1 provides a foundation for future work on scalable and unified RL training for agentic LLMs.

References

- [1] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai:

- Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
 - [3] Mingyue Cheng, Yucong Luo, Jie Ouyang, Qi Liu, Huijie Liu, Li Li, Shuo Yu, Bohou Zhang, Jiawei Cao, Jie Ma, et al. A survey on knowledge-oriented retrieval-augmented generation. *arXiv preprint arXiv:2503.10677*, 2025.
 - [4] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
 - [5] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
 - [6] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shiron Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
 - [7] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024.
 - [8] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*, 2020.
 - [9] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2023.
 - [10] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022.
 - [11] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
 - [12] Chuang Jiang, Mingyue Cheng, Xiaoyu Tao, Qingyang Mao, Jie Ouyang, and Qi Liu. Tablemind: An autonomous programmatic agent for tool-augmented table reasoning. *arXiv preprint arXiv:2509.06278*, 2025.
 - [13] Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.
 - [14] Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023.
 - [15] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
 - [16] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.

- [17] Yucong Luo, Yitong Zhou, Mingyue Cheng, Jiahao Wang, Daoyu Wang, Tingyue Pan, and Jintao Zhang. Time series forecasting as reasoning: A slow-thinking approach with reinforced llms. *arXiv preprint arXiv:2506.10630*, 2025.
- [18] Jie Ouyang, Tingyue Pan, Mingyue Cheng, Ruiran Yan, Yucong Luo, Jiaying Lin, and Qi Liu. Hoh: A dynamic benchmark for evaluating the impact of outdated information on retrieval-augmented generation. *arXiv preprint arXiv:2503.04800*, 2025.
- [19] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [20] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22, 2023.
- [21] Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, et al. Kilt: a benchmark for knowledge intensive language tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544, 2021.
- [22] Jiahao Qiu, Xuan Qi, Tongcheng Zhang, Xinzhe Juan, Jiacheng Guo, Yifu Lu, Yimin Wang, Zixin Yao, Qihan Ren, Xun Jiang, et al. Alita: Generalist agent enabling scalable agentic reasoning with minimal predefinition and maximal self-evolution. *arXiv preprint arXiv:2505.20286*, 2025.
- [23] Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, et al. Qwen2.5 technical report, 2025.
- [24] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [25] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
- [26] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [27] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- [28] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554, 2022.
- [29] Daoyu Wang, Mingyue Cheng, Qi Liu, Shuo Yu, Zirui Liu, and Ze Guo. Paperarena: An evaluation benchmark for tool-augmented agentic reasoning on scientific literature. *arXiv preprint arXiv:2510.10909*, 2025.
- [30] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- [31] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345, 2024.

- [32] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.
- [33] Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Xin Guo, Dingwen Yang, Chenyang Liao, Wei He, et al. Agentgym: Evaluating and training large language model-based agents across diverse environments. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 27914–27961, 2025.
- [34] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 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, 2018.
- [35] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*, 2022.
- [36] Shuo Yu, Mingyue Cheng, Qi Liu, Daoyu Wang, Jiqian Yang, Jie Ouyang, Yucong Luo, Chenyi Lei, and Enhong Chen. Multi-source knowledge pruning for retrieval-augmented generation: A benchmark and empirical study. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management*, pages 3931–3941, 2025.
- [37] Shuo Yu, Mingyue Cheng, Daoyu Wang, Qi Liu, Zirui Liu, Ze Guo, and Xiaoyu Tao. Memweaver: A hierarchical memory from textual interactive behaviors for personalized generation. *arXiv preprint arXiv:2510.07713*, 2025.