

# ToolOmni: Enabling Open-World Tool Use via Agentic learning with Proactive Retrieval and Grounded Execution

Shouzheng Huang, Meishan Zhang, Baotian Hu<sup>✉</sup>, Min Zhang

Harbin Institute of Technology (Shenzhen)

huangshouzheng@stu.hit.edu.cn, mason.zms@gmail.com,

{hubaotian, zhangmin2021}@hit.edu.cn

<https://github.com/Huangsz2021/ToolOmni>

## Abstract

Large Language Models (LLMs) enhance their problem-solving capability by utilizing external tools. However, in **open-world** scenarios with massive and evolving tool repositories, existing methods relying on static embedding retrieval or parameter memorization of tools struggle to align user intent with tool semantics, leading to suboptimal accuracy of open-world tool retrieval and execution. To address these, we present **ToolOmni**, a unified agentic framework that enables LLMs for open-world tool use by proactive retrieval and grounded execution within a reasoning loop. First, we construct a cold-start multi-turn interaction dataset to instill foundational agentic capabilities via Supervised Fine-Tuning (SFT). Then, we introduce open-world tool learning based on a **Decoupled Multi-Objective GRPO** algorithm, which simultaneously optimizes LLMs for both tool retrieval accuracy and execution efficacy in online environments. Extensive experiments demonstrate that ToolOmni achieves state-of-the-art performance both in retrieval and execution, surpassing strong baselines by a significant margin of **+10.8%** in end-to-end execution success rate, while exhibiting exceptional robustness and generalization capabilities.

## 1 Introduction

Tool Learning with LLM achieves higher accuracy, efficiency, and autonomy in problem solving by combining the strengths of specialized tools and foundational models (Nakano et al., 2021; Yao et al., 2022, 2023; Schick et al., 2023; Qin et al., 2024; Wu et al., 2023; Li et al., 2025d). Efforts in this field have predominantly focused on teaching models to effectively use tools via demonstration-based learning, typically leveraging Supervised Fine-Tuning (SFT) on curated expert trajectories (Nakano et al., 2021; Yao et al., 2022; Schick et al.,

<sup>✉</sup>Corresponding author.

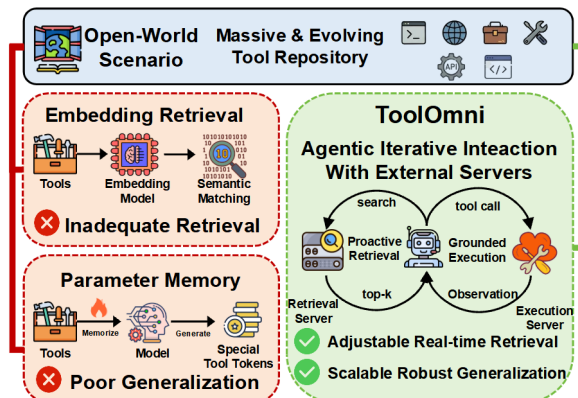


Figure 1: Motivation for **ToolOmni** in Open-World Scenarios: Embedding retrieval methods struggle with *Massive tools*, often resulting in low retrieval accuracy due to shallow matching; Parameter memory methods fail to adapt to *Evolving tools*, suffering from poor generalization to unseen tools. **ToolOmni** overcomes these limitations via a unified agentic framework that couples **Proactive Retrieval** with **Grounded Execution**, enabling effective open-world tool use.

2023; Li et al., 2025c), or feedback-based learning, which aligns model behavior with feedback from both environment and human through Reinforcement Learning (Schulman et al., 2017; Christiano et al., 2017; Nakano et al., 2021; Baker et al., 2022; Li et al., 2025b; Qian et al., 2025). However, in **open-world** scenarios characterized by massive and dynamically updated tool repositories, models must not only understand how to use tools but also master the ability to search and select the correct ones.

As shown in Fig.1, to help narrow the scope of relevant tools, prevailing solutions typically adopt a pipeline approach, employing embedding models to retrieve relevant tools based on query similarity before passing them to the execution agent (Qin et al., 2024; Chen et al., 2024; Xu et al., 2024). However, this paradigm operates passively, which decouples the retrieval process from the agent’s reasoning, preventing the execution model from

proactively participating in tool selection or refining the search based on task-specific needs. Consequently, these methods often struggle to effectively **align user intent** with the functionally essential tools in complex scenarios. Alternatively, some approaches (Wang et al., 2025; Schick et al., 2023; Su et al., 2025) fine-tune models to internalize tool documentation into parametric knowledge. While effective, this paradigm requires expensive retraining whenever the toolset updates, severely limiting **generalizability in dynamic environments**.

Recently, there has been a growing interest in agentic training frameworks that leverage Reinforcement Learning with Verifiable Rewards (RLVR) to unify reasoning with active environment interaction. Building on algorithms like Group Relative Policy Optimization (GRPO) (Shao et al., 2024) and Proximal Policy Optimization (PPO) (Schulman et al., 2017), recent works (Jin et al., 2025b; Xue et al., 2025; Qian et al., 2025; Li et al., 2025b) have demonstrated that LLMs can be trained to iteratively interact with the external environment, effectively invoking tools and leveraging feedback to enhance both performance and generalization. Nevertheless, these works often restrict LLMs to a limited toolset such as search engines and code executors, which constrains their applicability to diverse and dynamic open-world scenarios.

To address the above issues, we propose a unified agentic framework **ToolOmni** that integrates proactive retrieval and grounded execution into a unified end-to-end process, scaling agentic capabilities to dynamic, open-world tool scenarios. To build ToolOmni, we first conduct **Tool Learning cold start** phase with a high-quality hybrid dataset that integrates both retrieval and execution trajectories. This stage enables the model to acquire the foundational capabilities required for effective tool interaction. Building on this foundation, the second stage adopts an **Open World Tool Learning** process based on an enhanced Group Relative Policy Optimization algorithm. Unlike naive GRPO, which relies on a single reward, our framework treats retrieval and execution as interconnected yet distinct sub-tasks, where we compute task-specific rewards and advantages for retrieval and execution independently, and integrate them into a single optimization, enabling the synchronized optimization of both capabilities. This is crucial because decoupling allows us to provide finer-grained process supervision, ensuring that both retrieval recall

and execution reasoning are optimized precisely without mutual interference.

To evaluate the effectiveness of ToolOmni, we conduct extensive experiments on the ToolBench benchmark. The results demonstrate that ToolOmni achieves superior performance in both tool retrieval and task execution, particularly in open-world scenarios with massive candidate tools (+11.9%). Additionally, ToolOmni exhibits exceptional robustness to unseen instructions and tools, demonstrating that it learns universal tool-use mechanisms rather than relying on rigid memorization. The contributions of this work can be summarized as follows:

- We introduce **ToolOmni**, an end-to-end tool agentic framework that integrates proactive tool retrieval with grounded execution within a unified reasoning loop.
- We propose a two-stage training strategy that integrates a supervised cold-start for foundational tool retrieval and execution with GRPO-based RL for the subsequent synchronized optimization.
- Extensive experiments demonstrate that **ToolOmni** not only achieves superior performance in both tool retrieval and execution, but also shows strong robustness and generalizability to unseen domains.

## 2 Related Works

### 2.1 Tool Retrieval in Open-World Scenarios

In open-world scenarios, a common approach for tool retrieval is to employ an embedding model (Shi et al., 2025; Robertson et al., 2009; Qin et al., 2024; Qu et al., 2024; Zhao et al., 2025, 2026) that retrieves top-k relevant tools based on semantic similarity, narrowing the scope of candidate tools. Some works (Chen et al., 2024; Xu et al., 2024; Zheng et al., 2024; Kachuee et al., 2025) improve retrieval performance by using LLM to rewrite queries or expand tool documentation. Alternatively, another method trains the LLM to encode tool information into its parametric knowledge, enabling the model to directly generate corresponding tool identifiers to accomplish retrieval (Wang et al., 2025; Su et al., 2025). Our approach also utilizes embedding-based retrieval, yet transforms the interaction paradigm where the agent proactively formulates queries and invokes the embedding model as an executable tool.

## 2.2 LLM Tool Execution

Numerous studies have focused on augmenting LLMs with external tools to enhance their specialization and efficiency in solving complex tasks (Liang et al., 2024; Xu et al., 2023; Schick et al., 2023; Yao et al., 2022; Nakano et al., 2021). ReAct (Yao et al., 2023) establishes a Thought-Action-Observation loop, where the model generates explicit reasoning traces to justify and orchestrate tool execution in an interleaved manner. ToolLLM (Qin et al., 2024) adopts a tree search framework (DFSDT), allowing the model to explore multiple execution paths and back-track based on tool feedback to solve multi-step tasks. Meta-Tool (Qin et al., 2025) incorporates a plug-and-play retrieval module allowing the model to proactively search for relevant APIs, yet it treats retrieval and execution as isolated pipeline stages without joint optimization. More recently, agentic training frameworks that leverage RLVR have been employed to enhance the LLM’s ability to use external tools (Jin et al., 2025b; Xue et al., 2025; Qian et al., 2025; Li et al., 2025b,a). These studies treat external search and code execution as executable tools, facilitating an agentic, end-to-end reasoning paradigm for complex task solving. Inspired by them, ToolOmni broadens this scope to open-world by unifying proactive tool discovery and execution into an end-to-end process for complex task resolution.

## 3 Methodology

### 3.1 Problem Formulation

The goal of open-world tool-use agents is to address a user query  $Q$  by interacting with a large-scale tool repository  $\mathcal{T} = \{t_1, t_2, \dots, t_N\}$ , where  $N$  is the number of candidate tools. Given the massive scale of  $\mathcal{T}$  (e.g.,  $N > 10,000$ ), ToolOmni models the open-world tool-use process as a cascaded retrieval-execution framework driven by the policy  $\pi_\theta$ . At each turn  $t$ , the agent first performs an iterative *Proactive Retrieval* phase to identify a task-complete set of tools, followed by a *Grounded Execution* phase to invoke them. The trajectory sequence  $\tau$  and final answer  $a$  can be formulated as follows:

$$\tau = \left\{ \underbrace{(r_t^{ret}, \alpha_t^{ret}, \mathcal{T}_{sub})}_{\text{Proactive Retrieval}}, \underbrace{(r_t^{exe}, \alpha_t^{exe}, o_t)}_{\text{Grounded Execution}} \right\}_{t=1}^T, \\ a \sim \pi_\theta(Q, \tau) \quad (1)$$

where the cascaded retrieval-execution framework at each turn is determined by the policy:

$$\tau = \begin{cases} (r_t^{ret}, \alpha_t^{ret}) \sim \pi_\theta(\cdot | s_t, Q), & \text{Ret.} \\ (r_t^{exe}, \alpha_t^{exe}) \sim \pi_\theta(\cdot | s_t, Q, \mathcal{T}_{sub}), & \text{Exec.} \end{cases}$$

At each turn  $t$ , the agent’s state  $s_t$  consists of the history of all previous actions and the corresponding observations, i.e.,  $s_t = (a_1, s_1, \dots, a_{t-1}, o_{t-1})$ . As the Fig.2 shows, the agent iteratively interleaves reasoning  $r_i^{ret}$  and actions  $\alpha_i^{ret}$  in the Proactive Retrieval phase, where each action  $\alpha_i^{ret}$  either queries  $q_i$  to retrieve relevant tools or finalizes the task-complete toolset  $\mathcal{T}_{sub}$ . After finalizing the toolset  $\mathcal{T}_{sub}$ , the agent enters the Grounded Execution phase, where it interleaves reasoning  $r_t^{exe}$  and tool-call actions  $\alpha_t^{exe}$  grounded in the retrieved tool documentation to generate observations  $o_t$  or the final answer  $a$ .

### 3.2 Cold-start Tool Learning

To endow ToolOmni with basic tool-use capabilities, we first perform a SFT phase, serving as a cold-start initialization for the model.

**Data Curation.** Our training dataset is derived from ToolBench (Qin et al., 2024) and carefully curated to support both retrieval and execution phases. For tool retrieval, we first select a subset of 80,000 queries to train a specialized retrieval-only model. Using this model, we perform rejection sampling to remove low-quality instances, resulting in a high-quality corpus of approximately 28,000 retrieval trajectories. For tool execution, we extract around 33,000 trajectories from the ToolBench training set, including both correct and incorrect execution paths. To ensure data quality, we employ Qwen-2.5-32B as an automated judge to rigorously validate each trajectory. Further details are presented in Appendix A.

**SFT Objective.** Given the dataset of high-quality trajectories  $\mathcal{D} = \{(\tau_{ret}, \tau_{exe})\}$ , we optimize the model using the standard cross-entropy loss:

$$\mathcal{L}_{SFT} = - \sum_{(x,y) \in \mathcal{D}} \log \pi_\theta(y | x) \quad (2)$$

Following prior works (Jin et al., 2025b; Huang et al., 2025; Jin et al., 2025a; Sun et al., 2025), we compute token-level losses exclusively on the agent-generated reasoning traces  $r$  and actions  $\alpha$ , masking all external observations  $o$ . This strategy prevents the model from predicting environmental

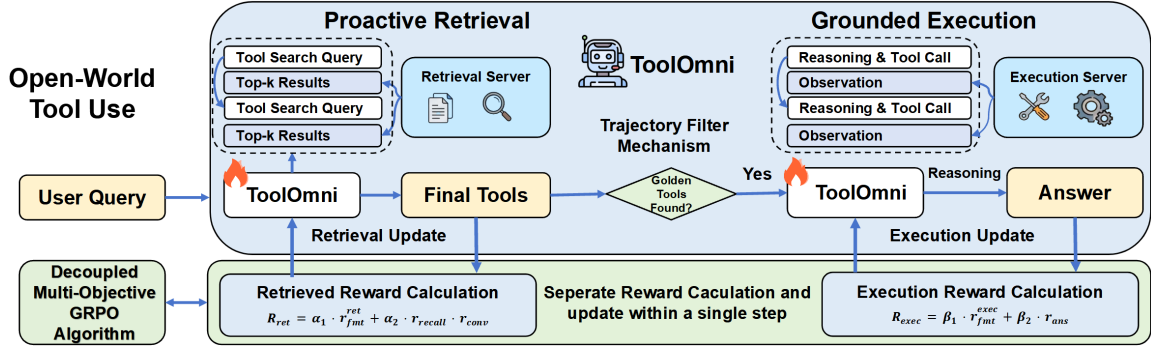


Figure 2: Overview of the **ToolOmni** framework. The pipeline operates in two decoupled phases: **Proactive Retrieval**: The agent iteratively interacts with the retrieval server to curate a candidate tool set. **Grounded Execution**: With retrieval results, the agent performs reasoning and tool invocation to generate the final answer.

dynamics, thereby stabilizing the policy gradient optimization.

### 3.3 Open-world Tool Learning

While SFT establishes basic tool-use capabilities, it relies on imitating high-quality labeled trajectories, restricting exploration of diverse trajectories and limiting scalability to large-scale, open-world tool scenarios. To bridge this gap, we introduce open-world tool learning based on GRPO, which enables the agent to optimize its proactive retrieval and grounded execution actions through iterative trial and error.

#### 3.3.1 Proactive Tool Retrieval

The agentic interaction begins with the proactive tool retrieval phase. Unlike existing tool learning approaches that rely on a single-turn passive retrieval, ToolOmni autonomously determines whether and what to retrieve. Formally, given the user instruction  $Q$ , the policy  $\pi_\theta$  analyzes the user’s intents and then formulates a tool search query  $q_t$  encapsulated within special tags: `<search> q_t </search>`. Upon receiving  $q_t$ , the retrieval server encodes it using a pre-trained embedding model  $E(\cdot)$  and then performs retrieval. It computes the cosine similarity between the query vector  $E_{q_t}$  and the pre-indexed tool embeddings, returning the top- $k$  candidates:

$$\mathcal{T}_{ret}^t = \text{top-k}(\cos(E_\tau, E_{q_t})), \quad (3)$$

where  $\cos(\cdot)$  denotes cosine similarity function. ToolOmni iterates proactive retrieval, autonomously generating multiple search queries as needed with real-time retrieval. Once a complete set of tools sufficient for the task has been collected, ToolOmni selects useful tools from the candidates

and arranges them into a sorted subset that is finalized within `<tool_call>` and `</tool_call>` tags.

$$\mathcal{T}_{sub} = \pi_\theta(P_{ret}, Q, \bigcup_t \mathcal{T}_{ret}^t), \quad (4)$$

#### 3.3.2 Grounded Tool Execution

After obtaining the sorted tool subset  $\mathcal{T}_{sub}$ , ToolOmni is instructed via the execution prompt  $P_{exec}$  to enter the execution phase. We employ a strategic trajectory filtering procedure to ensure both the training stability of the execution policy and the solvability of the queries. Specifically, we only retain trajectories where the generated subset  $\mathcal{T}_{sub}$  successfully recall all ground-truth tools ( $\mathcal{T}_{gold} \subseteq \mathcal{T}_{sub}$ ). Based on these valid retrieval results, the policy  $\pi_\theta$  then conducts multi-turn interleaved reasoning and tool invocation to solve the user’s query.

Specifically, at each step  $t$ , ToolOmni first conducts reasoning inside `<reasoning>` and `</reasoning>` tags. Subsequently, it invokes a tool by explicitly specifying the target function name and its required arguments, e.g., `<tool_call> {"tool_name": "genderize", "tool_input": {"name": "john"}} </tool_call>`. To ensure stable and efficient online reinforcement learning, we deploy an *LLM-based Tool Simulator* that emulates the environment’s feedback (Guo et al., 2024). The simulator produces realistic execution results for tool invocations, which are enclosed within `<information>` and `</information>` tags and appended to the ongoing context. The iterative cycle of reasoning and tool invocation continues until it derives the final solution, which is ultimately presented within `<answer>` and `</answer>` tags. Finally, the grounded tool execution process can be formulated as:

$$y = \pi_\theta(P_{exec}, Q, \mathcal{T}_{sub}, O_{tool}). \quad (5)$$

Given the open-ended nature of tool execution, rigid rule-based matching is insufficient to verify the final answer. Instead, we utilize a reward model to evaluate each trajectory holistically. A positive reward is assigned when the agent successfully invokes the required tools and produces a correct final answer.

### 3.3.3 Open-world Tool Use Rewards

In the RL process, the reward functions are utilized to steer the model towards desired properties. Below, we design retrieval and execution rewards to induce the model to perform effective proactive retrieval and reliable grounded execution.

The retrieval reward  $R_{ret}$  comprises three weighted components: ensuring format correctness, maximizing recall of ground-truth tools, and promoting effective conversion of retrieved information:

$$R_{ret} = \alpha_1 \cdot r_{fmt}^{ret} + \alpha_2 \cdot r_{rec} \cdot r_{conv}. \quad (6)$$

Here,  $r_{fmt}^{ret} \in \{0, 1\}$  ensures format compliance;  $r_{rec} \in [0, 1]$  measures the recall of ground-truth tools  $\mathcal{T}_{gold}$  within the cumulative retrieved set  $\mathcal{T}_{ret}$ ; and  $r_{conv} \in [0, 1]$  rewards the proportion of these retrieved tools that are ultimately selected into the final set, encouraging effective utilization.

The execution reward  $R_{exec}$  is similarly formulated as a weighted combination of format and outcome correctness:

$$R_{exec} = \beta_1 \cdot r_{fmt}^{exec} + \beta_2 \cdot r_{ans}. \quad (7)$$

The term  $r_{fmt}^{exec} \in \{0, 1\}$  enforces structural validity of the reasoning path and tool invocations, while  $r_{ans} \in \{0, 1\}$  is a binary signal from the reward model indicating the correctness of the answer.

### 3.3.4 Online Policy Optimization

**Group-Relative Advantage Estimation.** For each query  $x$ , the agent generates a group of  $G$  outputs  $\{o_1, \dots, o_G\}$  sampled from the current policy  $\pi_\theta$ . We compute the advantages for retrieval and execution independently using the group-relative standard:

$$A_{task}^{(i)} = \frac{R_{task}^{(i)} - \text{mean}(\{R_{task}\}_G)}{\text{std}(\{R_{task}\}_G) + \epsilon} \quad (8)$$

where  $task \in \{ret, exec\}$ . By normalizing rewards within the group specifically for each sub-task, we isolate the learning signals, preventing the sparsity of execution rewards from destabilizing the retrieval learning process.

**Decoupled Policy Gradient.** Based on the estimated advantages, we optimize the policy  $\pi_\theta$  by maximizing the surrogate objective. Consistent with the SFT stage, we apply token-level masking to ensure that gradients are back-propagated solely through the agent-generated tokens. The final optimization objective is formulated as:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \mathbb{I}(t \in \mathcal{M}_{task}) \min \left( \rho_{i,t}(\theta) \hat{A}_i, \text{clip}(\rho_{i,t}(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right) \right]$$

**Optimization Stability.** To further stabilize the training, we implement two key optimizations for GRPO in this context: **(1) Separated Update:** Instead of summing the gradients directly, we perform the updates for retrieval and execution sequentially within a single step. This prevents gradient conflict where the magnitude of one objective overwhelms the other; **(2) Selective Rollout:** As detailed in Sec. 3.3.2, we enforce a strict filtering mechanism where the execution generation is initiated only when the retrieval stage successfully recalls all golden tools ( $\mathcal{T}_{gold} \subseteq \mathcal{T}_{ret}$ ). By excluding invalid retrieval instances prior to rollout, we ensure that the execution policy is trained exclusively on high-quality, grounded contexts.

## 4 Experiment

In this section, we conduct extensive experiments on the ToolBench benchmark to answer the following Research Questions (RQs): **RQ1:** How does ToolOmni perform in retrieving relevant tools from massive, open-world repositories compared to existing baselines? **RQ2:** Can ToolOmni effectively execute complex, multi-step tasks in an end-to-end manner, surpassing pipeline and unified approaches? **RQ3:** Does the proposed framework demonstrate robustness against retrieval noise and generalization capabilities across unseen tools and domains? **RQ4:** What are the individual contributions of the core components (e.g., iterative retrieval, training stages) and hyperparameters to the overall performance?

### 4.1 Experimental Setup

We evaluate ToolOmni on **ToolBench** (Qin et al., 2024; Guo et al., 2024), testing across three difficulty levels (I1–I3) and three generalization splits (Instruction, Tool, Category) to assess robustness.

Method	I1			I2			I3			Avg.
	@1	@3	@5	@1	@3	@5	@1	@3	@5	
<i>In-Domain</i>										
BM25*	29.46	31.12	33.27	24.13	25.29	27.65	32.00	25.88	29.78	28.73
EmbSim*	63.67	61.03	65.37	49.11	42.27	46.56	53.00	46.40	52.73	53.35
Re-Invoke*	69.47	-	61.10	54.56	-	53.79	59.65	-	59.55	59.69
IterFeedback	71.64	71.29	76.31	62.65	55.58	60.62	73.85	64.06	69.02	67.22
ToolGen	69.47	72.26	79.12	46.77	53.58	62.45	77.06	<u>76.44</u>	<b>85.48</b>	69.18
ToolRetriever	<u>81.91</u>	<u>82.05</u>	<b>85.57</b>	<u>75.92</u>	<u>69.61</u>	<b>75.40</b>	<u>79.82</u>	72.42	77.11	<u>77.76</u>
<b>ToolOmni</b>	<b>86.07</b>	<b>84.49</b>	<u>85.56</u>	<b>81.50</b>	<b>73.13</b>	<u>74.11</u>	<b>84.40</b>	<b>76.81</b>	<u>77.25</u>	<b>80.37</b>
<i>Multi-Domain</i>										
BM25*	22.77	22.64	25.61	18.29	20.74	22.18	10.00	10.08	12.33	18.29
EmbSim*	54.00	50.82	55.86	40.84	36.67	39.55	18.00	17.77	20.70	37.13
ToolGen	68.84	72.00	78.76	46.77	53.55	62.43	75.68	<b>75.26</b>	<b>84.48</b>	68.64
IterFeedback	71.77	70.49	74.82	63.18	55.41	61.32	67.43	59.23	63.14	65.20
ToolRetriever	<u>81.03</u>	<u>80.88</u>	<u>84.52</u>	<u>75.91</u>	<u>69.53</u>	<b>75.21</b>	<u>77.52</u>	69.15	<u>74.24</u>	<u>76.44</u>
<b>ToolOmni</b>	<b>86.07</b>	<b>84.09</b>	<b>84.80</b>	<b>80.63</b>	<b>72.96</b>	<u>73.91</u>	<b>78.90</b>	<u>71.44</u>	71.82	<b>78.29</b>

Table 1: Main results of tool retrieval performance (NDCG@ $k$  %) on ToolBench. **ToolOmni** achieves the best performance across most metrics. The best performance is boldfaced, while the second-best performance is underlined. **Note:** Methods marked with \* report results cited from ToolGen (Wang et al., 2025).

Retrieval performance is measured by **NDCG@ $k$**  ( $k \in \{1, 3, 5\}$ ), while end-to-end execution is evaluated using **Solvable Pass Rate (SoPR)** and **Win Rate (SoWR)**, computed by a GPT-5 judge. We initialize our model with **Qwen3-4B-Instruct** (Yang et al., 2025) and train it via the proposed decoupled GRPO ( $G = 5, T = 1.0$ ) on 8 NVIDIA H100 GPUs. For comparison, we benchmark against competitive baselines across both tasks. **For retrieval**, we consider sparse (BM25 (Robertson et al., 2009)), dense (ToolRetriever (Qin et al., 2024), EmbSim<sup>1</sup>), and refinement methods (IterFeedback (Xu et al., 2024), Re-Invoke (Chen et al., 2024)), alongside the generative ToolGen (Wang et al., 2025). **For execution**, we evaluate pipeline agents (GPT, ToolLLaMA (Qin et al., 2024)) paired with ToolRetriever, and the unified generative model ToolGen (Wang et al., 2025). Further details are provided in **Appendix B**.

## 4.2 Tool Retrieval Performance(RQ1)

We evaluate retrieval under two settings: *In-Domain*, where the search space is restricted to the specific subset of tools relevant to the current test group; and *Multi-Domain*, where the agent must identify the correct tools from the entire repository of over **16,000** APIs. Regarding the mechanism, baselines typically employ static one-shot retrieval. ToolGen uses generative retrieval, while IterFeedback performs multi-turn query refinement. Simi-

<sup>1</sup>OpenAI’s sentence embedding model:text-embedding-3-large

larly, ToolOmni adopts an agentic iterative retrieval paradigm. To ensure a fair comparison, we align the search budget of ToolOmni with IterFeedback, limiting both to a maximum of 4 retrieval turns.

As shown in Table 1, **ToolOmni** achieves superior performance across the majority of metrics, particularly in the most rigorous Multi-Domain setting where it attains the highest average NDCG of **78.29%**. It significantly outperforms strong baselines like ToolRetriever and ToolGen in top-1 (@1) and top-3 (@3) precision, demonstrating its superior ability to accurately pinpoint the "golden tools" from a massive repository.

Notably, ToolOmni occasionally scores slightly lower on **NDCG@5**. This is an expected side effect of our *proactive selection mechanism*. Unlike standard retrievers that always return a fixed top- $k$  list, ToolOmni autonomously decides when to stop searching. If a task requires only 1 or 2 tools, our model outputs a concise set rather than padding it with irrelevant candidates to fill the top-5 slots. Thus, this performance gap reflects a preference for efficiency and precision over merely maximizing recall metrics.

## 4.3 Tool Execution Performance(RQ2)

We assess execution performance in two scenarios: **(1)With Golden Truth** and **(2)With Retriever** (end-to-end setting). For the end-to-end comparison, pipeline baselines are paired with ToolRetriever to perform a single initial search. ToolOmni and ToolGen utilize their internal retrieval mecha-

Method	Solvable Pass Rate (SoPR)				Solvable Win Rate (SoWR)			
	I1	I2	I3	Avg.	I1	I2	I3	Avg.
<i>With Ground-Truth Tools</i>								
GPT-3.5	56.60	47.80	<b>54.64</b>	53.01	-	-	-	-
ToolLlama-v2	43.60	45.80	39.30	42.90	38.65	55.66	32.79	42.37
ToolLlama-v1	24.20	29.20	12.30	21.90	28.22	42.45	16.39	29.02
ToolGen (3B)	40.80	43.40	21.30	35.17	39.26	47.17	14.75	33.73
ToolGen (7B)	44.20	42.90	32.00	39.70	35.58	45.28	16.39	32.42
<b>ToolOmni</b>	<b>60.70</b>	<b>50.90</b>	<u>52.50</u>	<b>54.70</b>	<b>52.76</b>	<b>56.60</b>	<b>40.98</b>	<b>50.11</b>
<i>With Retriever (End-to-End)</i>								
GPT-3.5	<u>51.43</u>	41.19	34.43	42.35	43.56	46.23	29.51	39.77
ToolLlama-v2	43.10	<u>46.70</u>	<u>40.20</u>	<u>43.33</u>	42.94	<u>50.00</u>	14.75	35.90
ToolLlama-v1	28.05	29.20	13.90	23.72	26.38	33.96	11.48	23.94
ToolGen (3B)	37.40	40.60	30.30	36.10	36.20	29.25	9.84	25.10
ToolGen (7B)	39.60	40.60	23.00	34.40	30.06	29.25	16.39	25.23
<b>ToolOmni</b>	<b>59.20</b>	<b>58.10</b>	<b>45.10</b>	<b>54.13</b> $\uparrow$ <b>10.8</b>	<b>50.31</b>	<b>57.55</b>	<b>42.62</b>	<b>50.16</b> $\uparrow$ <b>10.5</b>

Table 2: End-to-end execution performance on ToolBench. We report both **Solvable Pass Rate (SoPR %)** and **Solvable Win Rate (SoWR %)**. ToolOmni demonstrates robust superiority, significantly outperforming GPT-3.5 and ToolLlama-v2 in end-to-end settings.

nisms. To ensure fair comparison during execution, we enforce a uniform interaction budget: the maximum number of tool execution turns is capped at 6 for all models.

Table 2 details the execution results. With Golden Tools, **ToolOmni** attains the highest average SoPR of **54.70%**, surpassing GPT-3.5 (53.01%) and ToolLlama-v2 (42.90%), validating its superior reasoning capabilities. In the End-to-End setting, this advantage widens significantly. ToolOmni achieves **54.13%** SoPR, outperforming the GPT-3.5 pipeline (42.35%) by **+11.78%** and doubling the gain over the unified model ToolGen (36.10%). These results confirm the efficacy of our decoupled optimization strategy.

Furthermore, **ToolOmni** demonstrates robust superiority regarding the response quality metric **Solvable Win Rate (SoWR)**. In the end-to-end setting, it achieves an average SoWR of **50.16%**, significantly surpassing the strongest open-source baseline ToolLlama-v2 (35.90%) and the GPT-3.5 pipeline (39.77%). This indicates that ToolOmni not only correctly solves more queries but also generates more precise and coherent answers.

#### 4.4 Robustness Analysis(RQ3)

**Robustness Analysis on Generalization Splits.** As presented in Table 3, **ToolOmni** demonstrates exceptional robustness against distribution shifts. In the **Tool Generalization** setting (unseen tools within known categories), it achieves a SoPR of **52.20%**, effectively adapting to new API signatures. More impressively, in the rigorous **Category**

Method	Tool Gen.	Category Gen.	Average
GPT-3.5	50.60	42.10	46.35
ToolLlama-v2	39.20	39.40	39.30
ToolGen (7B)	33.20	40.10	36.65
<b>ToolOmni</b>	<b>52.20</b>	<b>55.95</b>	<b>54.08</b>

Table 3: Generalization performance (SoPR %) in the **End-to-End** setting. We report performance on Tool Gen. and Category Gen.. **ToolOmni** achieves the best generalization robustness.

**Generalization** setting (entirely novel domains), **ToolOmni** attains a remarkable score of **55.95%**. This performance surpasses a strong pipeline baseline ChatGPT (42.10%) by a substantial margin of **+13.85%**. This result indicates that while generative baselines like ToolGen tend to overfit to the specific tool seen during training, ToolOmni successfully learns the *universal meta-skills* of tool usage—such as parameter inference and error recovery. Consequently, it can transfer its reasoning capabilities to completely new domains without requiring domain-specific fine-tuning.

**Robustness against Retrieval Noise.** To test resilience, we injected adversarial tools that were retrieved by a dense model but excluding the ground truth—into the context at levels  $N \in \{0, 5, 10, 15\}$ . As shown in Figure 4, **ToolLlama-v2** degrades monotonically (39.3%  $\rightarrow$  20.5%), revealing its susceptibility to semantic distraction. In contrast, **ToolOmni** exhibits an adaptive resilience pattern. After an initial decline, its accuracy significantly **rebounds** to **58.2%** at  $N = 15$ . This indicates that as the candidate pool expands, our agent adap-

Settings	Ret. NDCG	Exec. SoPR
<b>Full</b>	<b>78.3</b>	<b>52.5</b>
w/o SFT	77.8	43.4
w/o RL	49.3	23.8

Table 4: Ablation on training stages. **Full ToolOmni** integrated with SFT as well as RL achieves the optimal results.

tively identifies functionally similar *alternative tools* within the adversarial set, showcasing a flexible reasoning capability that leverages noise to mitigate single-tool failures, transcending rigid ground-truth matching (Cuconasu et al., 2024).

#### 4.5 Ablation Study (RQ4)

To validate the contribution of our core components, we conduct comprehensive ablation studies across three critical dimensions: the multi-stage training pipeline, the iterative retrieval mechanism, and the specific reinforcement learning designs. For the macroscopic architecture, we first compare ToolOmni against three variants: (1) **w/o Iterative Retrieval**, a static baseline restricted to single-round retrieval; (2) **w/o RL**, the model trained solely via SFT; and (3) **w/o SFT**, the model trained directly via RL from the base LLM, skipping the cold-start phase. Furthermore, to disentangle the effectiveness of our Decoupled Multi-Objective GRPO algorithm, we introduce specific micro-ablations for the RL components, including **w/o Filter**, **Combined Update**, and **Vanilla GRPO**.

**Impact of Training Stages.** Table 4 disentangles the impact of training stages. For **Retrieval**, the removal of RL causes a drastic drop in NDCG (-28.99%), while removing SFT has a negligible impact (-0.51%), indicating that retrieval quality is primarily driven by RL-based exploration. Conversely, for **Execution**, the *w/o SFT* variant suffers a significant performance gap compared to the full model (43.40% vs. 52.50%). This confirms that SFT is essential as a cold-start mechanism, providing the necessary reasoning foundation that enables RL to optimize for complex tool-use scenarios effectively.

**Impact of Iterative Retrieval.** We further analyze the effectiveness of the retrieval strategy by comparing the Average NDCG@5 across all splits.

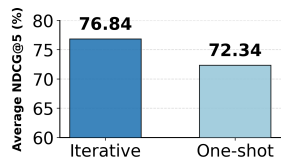


Figure 3: Ablation on retrieval strategy. **Iterative Retrieval** boosts NDCG@5 over the One-shot baseline.

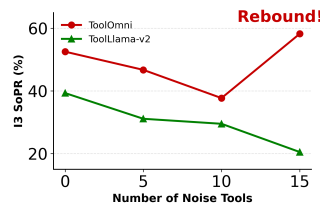


Figure 4: Robustness against adversarial noise on complex I3 tasks.

Top-k	Avg. NDCG
1	62.8
3	75.8
<b>5</b>	<b>78.3</b>
7	77.6
9	73.2

Table 5: Ablation study on retrieval strategy. “5” is the default setting.

As shown in Figure 3, the **Iterative Retrieval** mechanism employed by ToolOmni achieves an average score of **76.84%**, surpassing the static **One-shot Retrieval** baseline by **+4.5%**. This indicates that the ability to iteratively refine search queries allows the agent to dynamically adjust its search intent and filter out noise, effectively pinpointing the most utility-oriented tools required for execution.

**Impact of RL Components Design.** To validate the specific contributions of our proposed Decoupled Multi-Objective GRPO algorithm, we conduct a detailed ablation study on its key components.

As illustrated in Figure 5, our full **ToolOmni** model achieves the highest SoPR of 52.5%. The most substantial performance drop occurs when removing the trajectory filtering mechanism, falling to 38.5%, confirming that forcing the execution model to train on invalid contexts severely hinders its ability to reason and answer questions correctly. Furthermore, the **Combined Update** variant suffers a notable decline (42.6%), highlighting that directly summing gradients from the retrieval and execution phases leads to destructive gradient conflicts. Finally, replacing our decoupled reward design with **Vanilla GRPO** degrades performance (50.8%), demonstrating that assigning fine-grained, independent rewards provides more precise supervision than a sparse, unified reward.

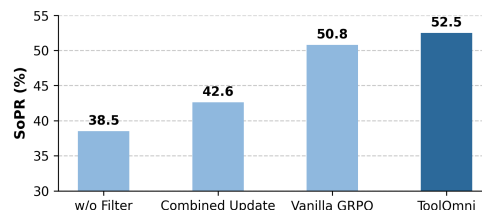


Figure 5: Ablation study of RL components on the ToolBench. Decoupled optimization and trajectory filtering are crucial for execution performance.

## 4.6 Hyperparameter Sensitivity(RQ4)

**Sensitivity to Retrieval Count ( $k$ ).** We investigate how the number of retrieved tools provided to the agent affects overall performance. Table 5 reports the Average NDCG across different retrieval counts  $k \in \{1, 3, 5, 7, 9\}$ .

The performance follows a clear **inverted U-shaped trajectory**. At low values ( $k = 1$ ), the model suffers significantly (62.84%), as the constrained context often fails to include the necessary ground-truth tools (low recall). As  $k$  increases to 5, performance peaks at **78.29%**, indicating an optimal balance where the context is rich enough to cover required functionalities without being overwhelming. However, increasing  $k$  further to 9 leads to a noticeable decline (73.16%). This degradation suggests that an excessively long context introduces irrelevant "noise tools," which dilutes the agent’s attention and complicates the reasoning process. Thus, our default setting of  $k = 5$  proves to be the most robust configuration for open-world scenarios.

**Sensitivity to Format Reward Weight.** We further investigate the impact of the format reward weight on execution performance. Specifically, We varied the format weight across  $\{0.2, 0.4, 0.6, 0.8\}$ . The results demonstrate a clear degradation trend when the format constraint becomes too dominant. While a moderate weight of 0.4 yields a peak SoPR of 44.3%, further increasing the weight to 0.6 and 0.8 causes the SoPR to drop significantly to 41.0% and 38.5%, respectively. This indicates that an excessively high format reward forces the model to overly prioritize strict syntactic compliance at the expense of exploring complex reasoning paths and functional tool chains. Consequently, a lower format weight strikes the optimal balance, ensuring structural validity without stifling the agent’s problem-solving capabilities in intricate open-world scenarios.

## 4.7 Computational Efficiency

A potential concern regarding the iterative proactive retrieval mechanism is the introduction of additional computational overhead. To address this, we analyze the overall efficiency by comparing the average number of search queries and actual tool execution calls per instruction across different models on the ToolBench benchmark.

As shown in Table 6, ToolOmni performs more search operations (3.02) than static baselines (1.00),

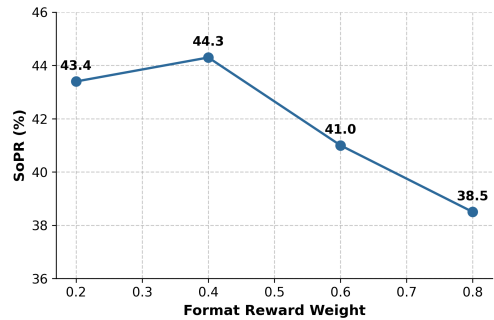


Figure 6: Sensitivity analysis of the format reward weight on the ToolBench.

Metric	ToolLlama	ToolGen	ChatGPT	ToolOmni
Search	1.00	-	1.00	3.02
Tool	2.40	2.75	3.98	2.47

Table 6: Average number of search times and tool execution calls per query on the ToolBench Benchmark.

However, this proactive exploration effectively filters noise, resulting in significantly fewer actual tool invocations (2.47 vs. ChatGPT’s 3.98). In real-world deployments, external API calls are the primary system bottleneck due to network latency and strict rate limits, whereas local embedding retrieval is computationally inexpensive. By trading a slight increase in cheap retrievals for a substantial reduction in costly, error-prone API calls, ToolOmni optimizes overall computational efficiency without sacrificing task success.

## 5 Conclusion

In this paper, we presented **ToolOmni**, a comprehensive agentic framework that enables open-world tool use via agentic learning with proactive retrieval and grounded execution. To address the inherent challenges of sparse rewards and sequential dependencies in this domain, we proposed a novel **Decoupled Multi-Objective GRPO** algorithm. Extensive experiments on the ToolBench benchmark demonstrate that **ToolOmni** achieves state-of-the-art performance in retrieval and end-to-end generation. Notably, our agent exhibits exceptional robustness against retrieval noise and strong generalization capabilities across unseen tools and domains. These results validate that equipping LLMs with the ability to iteratively reason about tool selection and execution is the key to scalable and robust tool learning. Future work will explore extending this decoupled framework to multi-modal tools and more complex, long-horizon planning tasks.

## Limitations

While ToolOmni demonstrates strong generalization, we acknowledge two primary limitations. First, our architecture relies on a rigid cascaded paradigm. While stable, this separation restricts the agent’s flexibility in extremely complex tasks that require on-the-fly tool chain discovery based on intermediate execution results. Second, due to computational constraints, ToolOmni is currently trained exclusively on the Qwen3-4B base model. Consequently, the upper bound of its performance and the potential emergent reasoning capabilities from scaling our framework to larger foundation models remain unexplored.

## Acknowledgements

This work is jointly supported by National Natural Science Foundation of China (Grant No.62422603), the Guangdong Basic and Applied Basic Research Foundation (Grant No.2024B0101050003), and Shenzhen Science and Technology Program (Grant No.ZDSYS20230626091203008). We sincerely thank all anonymous reviewers and Area Chairs for their detailed and careful reviews and valuable suggestions, which have significantly improved our work.

## References

Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. 2022. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654.

Yanfei Chen, Jinsung Yoon, Devendra Singh Sachan, Qingze Wang, Vincent Cohen-Addad, Mohammad-Hossein Bateni, Chen-Yu Lee, and Tomas Pfister. 2024. Re-invoke: Tool invocation rewriting for zero-shot tool retrieval. In *EMNLP (Findings)*, pages 4705–4726. Association for Computational Linguistics.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.

Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonello, and Fabrizio Silvestri. 2024. The power of noise: Redefining retrieval for RAG systems. In *SIGIR*, pages 719–729. ACM.

Zhicheng Guo, Sijie Cheng, Hao Wang, Shihao Liang, Yujia Qin, Peng Li, Zhiyuan Liu, Maosong Sun, and

Yang Liu. 2024. Stabletoolbench: Towards stable large-scale benchmarking on tool learning of large language models. *arXiv preprint arXiv:2403.07714*.

Ziyang Huang, Xiaowei Yuan, Yiming Ju, Jun Zhao, and Kang Liu. 2025. Reinforced internal-external knowledge synergistic reasoning for efficient adaptive search agent. *arXiv preprint arXiv:2505.07596*.

Bowen Jin, Jinsung Yoon, Priyanka Kargupta, Sercan O Arik, and Jiawei Han. 2025a. An empirical study on reinforcement learning for reasoning-search interleaved llm agents. *arXiv preprint arXiv:2505.15117*.

Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025b. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*.

Mohammad Kachuee, Sarthak Ahuja, Vaibhav Kumar, Puyang Xu, and Xiaohu Liu. 2025. Improving tool retrieval by leveraging large language models for query generation. In *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pages 29–38.

Maurice G Kendall. 1938. A new measure of rank correlation. *Biometrika*, 30(1-2):81–93.

Xiaoxi Li, Wenxiang Jiao, Jiarui Jin, Guanting Dong, Jijie Jin, Yinuo Wang, Hao Wang, Yutao Zhu, Ji-Rong Wen, Yuan Lu, and 1 others. 2025a. Deepagent: A general reasoning agent with scalable toolsets. *arXiv preprint arXiv:2510.21618*.

Xuefeng Li, Haoyang Zou, and Pengfei Liu. 2025b. Torl: Scaling tool-integrated rl. *arXiv preprint arXiv:2503.23383*.

Yunxin Li, Xinyu Chen, Shenyuan Jiang, Haoyuan Shi, Zhenyu Liu, Xuanyu Zhang, Nanhao Deng, Zhenran Xu, Yicheng Ma, Meishan Zhang, and 1 others. 2025c. Uni-moe-2.0-omni: Scaling language-centric omnimodal large model with advanced moe, training and data. *arXiv preprint arXiv:2511.12609*.

Yunxin Li, Zhenyu Liu, Zitao Li, Xuanyu Zhang, Zhenran Xu, Xinyu Chen, Haoyuan Shi, Shenyuan Jiang, Xintong Wang, Jifang Wang, and 1 others. 2025d. Perception, reason, think, and plan: A survey on large multimodal reasoning models. *arXiv preprint arXiv:2505.04921*.

Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, and 1 others. 2024. Taskmatrix. ai: Completing tasks by connecting foundation models with millions of apis. *Intelligent Computing*, 3:0063.

Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, and 1 others. 2021. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*.

- Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. 2025. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*.
- Shengqian Qin, Yakun Zhu, Linjie Mu, Shaoting Zhang, and Xiaofan Zhang. 2025. Meta-tool: Unleash open-world function calling capabilities of general-purpose large language models. In *ACL (1)*, pages 30653–30677. Association for Computational Linguistics.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. Toolllm: Facilitating large language models to master 16000+ real-world apis. In *ICLR*. OpenReview.net.
- Changle Qu, Sunhao Dai, Xiaoqi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2024. Towards completeness-oriented tool retrieval for large language models. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 1930–1940.
- Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. In *NeurIPS*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- 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*.
- Zhengliang Shi, Yuhan Wang, Lingyong Yan, Pengjie Ren, Shuaiqiang Wang, Dawei Yin, and Zhaochun Ren. 2025. Retrieval models aren’t tool-savvy: Benchmarking tool retrieval for large language models. *arXiv preprint arXiv:2503.01763*.
- Yunyue Su, Jingshuai Zhang, Bowen Fang, Wen Ye, Jinghao Zhang, Bowen Song, Weiqiang Wang, Qiang Liu, and Liang Wang. 2025. Toolscaler: Scalable generative tool calling via structure-aware semantic tokenization. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 556–578.
- Hao Sun, Zile Qiao, Jiayan Guo, Xuanbo Fan, Yingyan Hou, Yong Jiang, Pengjun Xie, Yan Zhang, Fei Huang, and Jingren Zhou. 2025. Zerosearch: Incentivize the search capability of llms without searching, 2025. URL <https://arxiv.org/abs/2505.04588>.
- Renxi Wang, Xudong Han, Lei Ji, Shu Wang, Timothy Baldwin, and Haonan Li. 2025. Toolgen: Unified tool retrieval and calling via generation. In *ICLR*. OpenReview.net.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*.
- Qiancheng Xu, Yongqi Li, Heming Xia, and Wenjie Li. 2024. Enhancing tool retrieval with iterative feedback from large language models. *arXiv preprint arXiv:2406.17465*.
- Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. 2023. On the tool manipulation capability of open-source large language models. *arXiv preprint arXiv:2305.16504*.
- Zhengkai Xue, Longtao Zheng, Qian Liu, Yingru Li, Xiaosen Zheng, Zejun Ma, and Bo An. 2025. Simpletir: End-to-end reinforcement learning for multi-turn tool-integrated reasoning. *arXiv preprint arXiv:2509.02479*.
- 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*.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In *ICLR*. OpenReview.net.
- Xinping Zhao, Xinshuo Hu, Zifei Shan, Shouzheng Huang, Yao Zhou, Xin Zhang, Zetian Sun, Zhenyu Liu, Dongfang Li, Xinyuan Wei, and 1 others. 2025. Kalm-embedding-v2: Superior training techniques and data inspire a versatile embedding model. *arXiv preprint arXiv:2506.20923*.
- Xinping Zhao, Xinshuo Hu, Jiaxin Xu, Danyu Tang, Xin Zhang, Mengjia Zhou, Yan Zhong, Yao Zhou, Zifei Shan, Meishan Zhang, and 1 others. 2026. Lmcb: Long-horizon memory embedding benchmark. *arXiv preprint arXiv:2603.12572*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623.

Yuanhang Zheng, Peng Li, Wei Liu, Yang Liu, Jian Luan, and Bin Wang. 2024. Toolrerank: Adaptive and hierarchy-aware reranking for tool retrieval. *arXiv preprint arXiv:2403.06551*.

## A Data Curation

### A.1 ToolBench

ToolBench (Qin et al., 2024) serves as a pioneering open-source benchmark for tool learning, constructed by scraping real-world APIs from RapidAPI. It encompasses a massive repository of **16,464 APIs** distributed across **49 distinct categories** (e.g., Social Media, E-commerce, Weather), reflecting a true open-world environment. Each API is documented with a JSON-based schema specifying its name, description, and calling parameters.

To generate high-quality instruction-tuning data, ToolBench utilizes a **Depth-First Search Decision Tree (DFS DT)** mechanism powered by ChatGPT. This process efficiently explores the vast action space to synthesize diverse and valid solution paths, resulting in over 126k tool-use instructions.

Crucially, these instructions cover a wide spectrum of complexity, ranging from single-tool tasks (I1) to intricate scenarios requiring the composition of multiple tools from the same category (I2) or across different collections (I3). This hierarchical diversity makes it an ideal testbed for evaluating an agent’s generalization and reasoning capabilities.

### A.2 Source and Difficulty Stratification

Our dataset is derived from the **ToolBench** corpus. To construct a balanced and challenging dataset, we first assess the difficulty of each query. We employ **ToolRetriever** to perform a preliminary dense retrieval, returning the top-5 tool candidates for each user query. An instance is classified based on recall success:

- **Easy Data:** All ground-truth "golden tools" are successfully retrieved within the top-5 candidates.
- **Hard Data:** At least one golden tool is missed by the retriever.

We then perform stratified sampling, selecting **60,000 hard instances** and **20,000 easy instances**, totaling **80,000 queries**. This set serves as the foundational pool for our subsequent processing and is also utilized as the prompt source for the Reinforcement Learning (RL) stage.

### A.3 Retrieval Training Data

From the initial 80,000 queries, we train a specialized retrieval-only model to generate candidate search trajectories. We then apply **Rejection Sampling** on these generations, filtering out trajectories where the model fails to locate the correct tools. This rigorous filtering yields a high-quality corpus of approximately **28,000 retrieval trajectories**, which is used for the retrieval component of the SFT stage.

### A.4 Execution Training Data

For the execution phase, we process the same pool of 80,000 instances to align with our agentic format. Since original ToolBench data lacks explicit reasoning steps, we employ **Qwen-2.5-72B-Instruct** to augment the data by generating detailed reasoning paths before each tool invocation. To ensure correctness, we use **Qwen-2.5-32B** as an automated judge to verify the execution results. To improve the model’s discrimination ability during training, we construct a mixed dataset comprising:

- **Positive Samples (70%):** Trajectories verified as correct.
- **Negative Samples (30%):** Trajectories containing errors (e.g., hallucinated parameters or wrong tool selection).

This process results in a final execution dataset of approximately **33,000 trajectories**.

### A.5 Final Dataset Composition

The **SFT Dataset** is the union of the 28,000 retrieval trajectories and 33,000 execution trajectories. The **RL Dataset** utilizes the initial pool of 80,000 queries as prompts to drive the online exploration and optimization process.

## B Experiment Setup

### B.1 Datasets and Metrics

Our experiments are conducted on ToolBench (Qin et al., 2024; Guo et al., 2024), a comprehensive real-world tool benchmark containing 16000+ real-world apis. Following the split of (Qin et al., 2024), the test queries are categorized into three complexity levels: **I1** (single-tool instructions), **I2** (intra-category multi-tool instructions), and **I3** (intra-collection multi-tool instructions). To rigorously evaluate robustness, these scenarios are further

stratified into three generalization splits: **Instruction (Inst)** generalization (unseen queries with seen tools), **Tool** generalization (unseen tools within known categories), and **Category (Cate)** generalization (unseen tools from entirely new domains).

We evaluate tool retrieval performance with **NDCG@k** (with  $k \in \{1, 3, 5\}$ ). For end-to-end evaluation, we report two key metrics: **Solvable Pass Rate (SoPR)**, which measures the percentage of queries successfully solved by the agent, and **Solvable Win Rate (SoWR)**, which indicates the proportion of answers that outperform those generated by a reference baseline. Both metrics are computed using GPT-5 Mini as an automated judge to ensure consistent and scalable evaluation.

### B.1.1 Retrieval Test Dataset

The retrieval evaluation is conducted on the official ToolBench test set. Table 7 details the distribution of test queries across the three instruction complexity levels.

Complexity Level	Number of Queries
I1 (Single-tool)	796
I2 (Intra-category)	573
I3 (Intra-collection)	218
<b>Total</b>	<b>1587</b>

Table 7: Statistics of the **ToolBench Retrieval Test Set**. The dataset is categorized by instruction complexity (I1: Single-tool, I2: Intra-category, I3: Intra-collection).

### B.1.2 Execution Test Dataset

For the end-to-end execution evaluation, we utilize the curated **Solvable Test Queries** from StableToolBench to ensure robust assessment. The test set comprises six distinct generalization subsets, covering different levels of difficulty (G1, G2, G3) and generalization types (Instruction, Category, Tool). Table 8 presents the detailed statistics.

### B.1.3 Evaluation Reliability

To address potential biases and stochastic instability associated with utilizing LLM-as-a-Judge for execution metrics, we conduct a multi-faceted reliability analysis to ensure the objectivity and consistency of our automated evaluation (Zheng et al., 2023).

**Human-Model Alignment.** We randomly sampled a subset of evaluation set and conducted a comprehensive human evaluation. Three human

Difficulty	Subset Name	Count
<b>G1</b>	Instruction Gen. (I1)	163
	Tool Gen. (I1)	158
	Category Gen. (I1)	153
<b>G2</b>	Instruction Gen. (I2)	106
	Category Gen. (I2)	124
<b>G3</b>	Instruction Gen. (I3)	61
<b>Total</b>		<b>765</b>

Table 8: Statistics of the **StableToolBench Solvable Test Set**. Queries are stratified by generalization difficulty (G1, G2, G3) and type.

annotators scored the trajectories using a three-point scheme: -1 (unsolved), 0 (unsure), and 1 (solved). We then calculated the Kendall tau correlation coefficient (Kendall, 1938) between the GPT-5 judgments and human annotations. The correlation result reached **0.847**, indicating a strong consistency between our automated evaluation and human assessment, which validates the reliability of our main experimental results.

**Cross-Model Consistency.** To further mitigate potential model-selection bias and ensure that the performance gains of ToolOmni are not backbone-specific, we evaluated our approach using three different strong models as judges: GPT-5, Gemini-3, and Qwen2.5-32B. As shown in Table 9, ToolOmni consistently achieves the state-of-the-art overall performance across all judge models. While the absolute SoPR values vary across different judges (e.g., Qwen tends to give higher scores overall), the relative ranking of the methods remains stable, further confirming the robustness of our framework.

Method	GPT-5	Gemini-3	Qwen-32B	Mean
ChatGPT	41.68	32.18	54.60	42.82
Llama-v2	42.22	30.28	53.00	41.83
Llama-v1	28.47	21.92	46.55	32.31
ToolGen-3B	35.62	21.68	36.62	31.31
ToolGen-7B	34.73	15.62	32.83	27.73
ToolOmni	54.08	31.95	79.15	<b>55.06</b>
<b>Average</b>	39.47	25.61	50.46	-

Table 9: Average SoPR (%) on ToolBench judged by different LLM backbones.

## B.2 Comparison Baselines.

To validate the effectiveness of ToolOmni, we benchmark it against a comprehensive set of competitive methods. For tool retrieval, we compare with traditional sparse retrievers like BM25

(Robertson et al., 2009) and dense retrievers such as ToolRetriever (Qin et al., 2024). Additionally, we evaluate advanced query refinement strategies, including Re-Invoke (Chen et al., 2024) and Iter-Feedback (Xu et al., 2024), as well as ToolGen (Wang et al., 2025), which adopts a generative paradigm for retrieval. For tool execution, our baselines encompass both proprietary and open-source models: we employ ChatGPT as strong zero-shot references, and compare against ToolLLaMA (Qin et al., 2024), the state-of-the-art SFT baseline on ToolBench. We also include ToolGen (Wang et al., 2025) again to assess the performance of unified generative frameworks in end-to-end scenarios.

### B.3 Implementation Details.

We initialize ToolOmni upon the Qwen3-4B-Instruct (Yang et al., 2025). Regarding the reward configuration, we set the format weight to 0.2 and the performance weight to 0.8 for both phases (i.e.,  $\alpha_1 = 0.2, \alpha_2 = 0.8$  for retrieval;  $\beta_1 = 0.2, \beta_2 = 0.8$  for execution) For the reinforcement learning stage, we configure the GRPO group size to  $G = 5$  and set the sampling temperature to 1.0 to encourage exploration during training. The retrieval module employs ToolRetriever as the dense embedding model. For the execution stage, we implement a hybrid environment that integrates real-world tools with MirrorAPI (Guo et al., 2024); this configuration ensures both the authenticity of tool interactions and the robustness All models are trained for a single epoch on 8 NVIDIA H100 GPUs.

### B.4 Execution Environment

To ensure both the authenticity of tool interactions and the robustness of the evaluation, we implement a hybrid execution environment that synergizes real-world API response logs with the MirrorAPI simulator (Guo et al., 2024).

Specifically, the environment operates on a fallback-based retrieval logic: upon a tool invocation, the system first queries a comprehensive repository of historical real-world API call records. If a matching execution trace is found, the authentic response is returned directly to ensure maximum fidelity. In cases where real-world records are unavailable, the system seamlessly transitions to MirrorAPI, which serves as a high-fidelity proxy. Crucially, MirrorAPI is not merely a deterministic mapping; it is pre-trained on an extensive corpus of real-world API interactions, allowing it to faithfully

replicate the stochasticity and error distributions inherent in practical environments. This enables the simulation of various failure modes, such as missing parameters, authentication errors, and service timeouts, thereby providing a rigorous testbed for the model’s error recovery capabilities. To illustrate the model’s interaction with this dynamic environment, we present a representative reasoning trajectory of ToolOmni in Box B.4.

#### ToolOmni Reasoning Trajectory under Execution Failures

```

1. Initial Tool Call (Failure)
<tool_call> {"category": "Advertising",
"tool_name": "Reqres", "api_name": "Users",
"tool_input": {}} </tool_call>
2. Environment Feedback (Simulated Error)
<information> {"error": "Missing input parameters.", "response": ""} </information>
3. Model Reasoning & Adaptive Adjustment
<reasoning> The initial call to retrieve user lists failed due to missing parameters. I will now switch to the "User by id" tool, providing the specific ID '1' to recover from this error. </reasoning>
4. Refined Tool Call (Success)
<tool_call> {"category": "Advertising",
"tool_name": "Reqres", "api_name": "User by id",
"tool_input": {"id": "1"}} </tool_call>
5. Final Environment Feedback
<information> {"error": "", "response": '{"user_id": 1, "name": "John Doe", "email": "john.doe@example.com"}'} </information>

```

This execution flow confirms that our hybrid environment not only provides realistic feedback but also effectively evaluates the model’s ability to detect, reason about, and recover from execution failures in complex, real-world scenarios.

## C Case Study

To provide a qualitative understanding of ToolOmni’s superiority, we present comprehensive case studies across both Open-Domain and Oracle settings. In the **Open-Domain scenario** (Fig.7), we demonstrate how ToolOmni’s proactive iterative retrieval effectively filters out noise and locates critical tools that pipeline baselines often miss, thereby preventing downstream hallucinations. In the **With Golden Truth setting** (Fig.8), we highlight the agent’s robust reasoning capabilities. Specifically, Case 2 illustrates how ToolOmni autonomously diagnoses API errors (e.g., missing tokens) and strategically pivots to alternative tools. Together, these cases validate that ToolOmni is not merely a tool invoker, but a resilient problem solver capable of navigating the complexities of

open-world environments.

## D Learning Algorithm of ToolOmni

Algorithm 1 outlines the complete training procedure. A critical feature of our approach is the **Filtered Rollout mechanism** (Line 6-9), which acts as a quality gate. By initiating execution rollouts only when the retrieval phase successfully recalls the golden tools, we ensure that the execution policy is trained exclusively on grounded, solvable contexts. This prevents the model from learning "hallucination shortcuts" to compensate for missing information. Furthermore, the **Separated Optimization** (Phase 3) allows the retrieval and execution modules to evolve at their own pace, guided by their respective specialized rewards, thereby stabilizing the overall learning dynamics in the complex open-world environment.

## E Prompts

To facilitate reproducibility, we provide the full system prompts utilized in our experiments. As illustrated in Figure 9, we design two distinct prompts tailored for the decoupled phases:

- The **Retrieval Prompt** (Left) instructs the agent to act as a "search copilot," employing an iterative loop of query generation and information synthesis to identify the optimal set of tools from the massive repository.
- The **Execution Prompt** (Right) guides the agent to function as an "AutoGPT," leveraging the retrieved tools within a grounded reasoning framework to solve the user's query step-by-step.

These prompts are used consistently across both the SFT data generation and the RL training stages.

---

**Algorithm 1** Decoupled Multi-Objective GRPO Algorithm

---

**Require:** Policy  $\pi_\theta$ ; Dataset  $\mathcal{D}$ ; Group size  $G$ ; Learning rate  $\eta$ .

**Ensure:** Optimized policy  $\pi_{\theta^*}$ .

```
1: for each training iteration do
2:   Sample a batch of queries  $\{x_i\}_{i=1}^B \sim \mathcal{D}$ .
3:   Phase 1: Group Rollout
4:   for each query  $x_i$  do
5:     for  $j = 1$  to  $G$  do
6:       Generate retrieval trajectory  $q_{i,j} \sim \pi_{\theta_{old}}(\cdot|x_i)$ .
7:       Retrieve tools  $\mathcal{T}_{i,j}$  using query  $q_{i,j}$ .
8:       if  $\mathcal{T}_{gold} \subseteq \mathcal{T}_{i,j}$  then ▷ Trajectory Filtering
9:         Generate execution trajectory  $e_{i,j} \sim \pi_{\theta_{old}}(\cdot|x_i, \mathcal{T}_{i,j})$ .
10:        end if
11:      end for
12:    end for
13:    Phase 2: Decoupled Reward Calculation and Advantage Estimation
14:    for each query  $x_i$  and trajectory  $q_{i,j}, e_{i,j}$  do
15:      Compute retrieval reward:  $R_{ret}^{i,j} = \alpha_1 r_{fmt}^{ret} + \alpha_2 r_{rec} \cdot r_{conv}$ .
16:      Estimate advantages for retrieval:  $A_{ret}^{i,j} = \frac{R_{ret}^{i,j} - \mu(R_{ret}^i)}{\sigma(R_{ret}^i) + \epsilon}$ .
17:      Compute execution reward:  $R_{exec}^{i,j} = \beta_1 r_{fmt}^{exec} + \beta_2 r_{ans}$ .
18:      Estimate advantages for execution:  $A_{exec}^{i,j} = \frac{R_{exec}^{i,j} - \mu(R_{exec}^i)}{\sigma(R_{exec}^i) + \epsilon}$ .
19:    end for
20:    Phase 3: Separated Optimization
21:    Retrieval Update  $\pi_\theta$  parameters:  $\theta \leftarrow \theta + \eta \nabla_\theta(\mathcal{J}_{ret})$ .
22:    Execution Update  $\pi_\theta$  parameters:  $\theta \leftarrow \theta + \eta \nabla_\theta(\mathcal{J}_{exec})$ .
23:  end for
24: return  $\pi_\theta$ 
```

---

## Case 1: Open-Domain Scenario

**Instruction:** I want to explore different genres of movies and series available on streaming platforms. Can you provide me with the basic details and streaming availability of a movie or series with the keyword 'zombie' and released between the years 2000 and 2025? Also, fetch the details of the actors and crew associated with the title with Watchmode ID 3173903.

### ✘ Baseline (ToolLlama + ToolRetriever):

[ToolRetriever Retrieval Result (Top-5)]

- [ID:781] OTT - Search (Noise)
- [ID:782] OTT - Title Details
- [ID:4580] Watchmode - Cast & Crew
- [ID:778] OTT - Advanced Search (Wrong Tool for Streaming)
- [ID:13631] Watchmode - Person (Noise)

**Missing Golden Tool:** [ID:7585] Streaming Availability - Get Basic

[Execution Phase]

**Action 1:** OTT-Search(title="zombie", page=1)

**Error:** Used basic search; returned random titles like "Zombie Outbreak" without checking year/genre.

**Action 2:** TitleDetails(imdbid="tt13411672")

**Result:** Streaming information.

**Action 3:** TitleDetails(id="3173903")

**Error:** Wrong tool! Used 'TitleDetails' instead of 'CastCrew' for Watchmode ID. Returned unrelated anime "Jujutsu Kaisen".

**Outcome: Hallucination.** Falsely claims "Available on Netflix" and provides wrong cast info.

### 🏠 ToolOmni (Ours):

[Phase 1: Proactive Retrieval]

<search> fetch basic details and streaming availability of a movie or series </search>

<information> [ID:782] OTT - Title Details Get basic information ... </information>

<search> confirm the keyword 'zombie' and released between the years 2000 and 2025, specifically for streaming platforms. </search>

<information> [ID:7585] Streaming Availability - Get Basic Get details of a show ... </information>

<search> associated with Watchmode ID 3173903, specifically for streaming platforms and crew details. </search>

<information> [ID:4580] Watchmode - Title Cast & Crew Return all people associated with the title ... </information>

<final\_tools> ID: [782,7585,4580] </final\_tools>

[Phase 2: Grounded Execution]

<reasoning> To fulfill this request, I need to:1. Search for movies or series with the keyword 'zombie' and released between 2000 and 2025 using the OTT details search tool with a filter for release year.2. Fetch the basic details and streaming availability of such a movie or series using the Streaming Availability tool with a country parameter (e.g., US).3. Retrieve the cast and crew details for the title with Watchmode ID 3173903 using the Watchmode tool. </reasoning>

<tool\_call> {"tool": "OTT details", "api": "Search", "input": {"title": "zombie", "start\_year": 2000, "end\_year": 2025}} </tool\_call>

<information> Found "Zombie (2000)", "Zombie Nightmare"... </information>

<reasoning> Check streaming availability. </reasoning>

<tool\_call> {"tool": "Streaming Availability", "api": "Get Basic", "input": {"country": "us", "imdb\_id": "tt0238784"}} </tool\_call>

<information> Streaming Not Available </information>

<reasoning> Get crew details </reasoning>

<tool\_call> {"tool": "Watchmode", "api": "Title Cast & Crew", "input": {"id": "3173903"}} </tool\_call>

<information> Crew details ... </information>

<reasoning> The search for movies or series with the keyword 'zombie' and released between 2000 and 2025 returned valid results, including both movies and TV movies. Additionally, the cast and crew details for the title with Watchmode ID 3173903 have been successfully retrieved. Although streaming availability was not found in the US, the user's request is satisfied with the exploration of genres and the cast/crew details for the specified title. </reasoning>

<answer> Provided movie list (2000-2025), streaming status and correct cast list. </answer>

**Outcome: Solved**

Figure 7: Qualitative comparison (Case Study). **Case 1** illustrates how pipeline baselines fail due to retrieval noise (selecting basic Search instead of Advanced Search) and missing tools (Streaming), leading to hallucinations. In contrast, **ToolOmni**'s proactive iterative retrieval precisely locates the golden toolset, enabling correct execution and grounded response generation. *Note: The iterative process shown is condensed for brevity.*

## Case 2: Complex Reasoning With Golden Tools

**Instruction:** I'm a movie enthusiast and I'm interested in watching critically acclaimed films. Can you suggest some award-winning movies with high IMDb ratings? Also, provide the streaming platforms where these movies are available in the US.

### ✘ Baseline (ToolLlama):

[Step 1]

Action: Streaming.Search(country="US", services="netflix,prime...")

**Error:** Invalid country code.

[Step 2]

Action: Streaming.Search(country="US", ...) (Retrying same args)

**Error:** Invalid country code.

[Step 3-5]

(Repeatedly calls the same tool with minor variations, ignoring the error message.)

**Outcome: Give Up.** The model entered a repetitive loop and failed to adjust its strategy, eventually terminating with give\_up\_and\_restart.

### 🛠️ ToolOmni (Ours):

**<reasoning>** The user wants to watch critically acclaimed films with high IMDb ratings and wants to know the streaming platforms where these movies are available in the US. To fulfill this request, I should start by searching for movies that are both critically acclaimed and have high IMDb ratings. **</reasoning>**

**<tool\_call>** {"tool": "Streaming Availability", "api": "Search Ultra", "input": {"country": "us", "min\_imdb\_rating": 80}} **</tool\_call>**

**<information>** Empty **</information>**

**<reasoning>** Search returned empty... parameters might be too restrictive. I should try fetching "New Titles" via Watchmode or check a specific known movie. **</reasoning>**

**<tool\_call>** {"tool": "Watchmode", "api": "New Titles", "input": {"types": "movie"}} **</tool\_call>**

**<information>** Empty **</information>**

**<reasoning>** Still no results. Instead of relying on general search, I should directly verify a known critically acclaimed film like "The Shawshank Redemption". **</reasoning>**

**<tool\_call>** {"tool": "OTT details", "api": "Title Details", "input": {"imdbid": "tt0111161"}} **</tool\_call>**

**<information>** Found "The Shawshank Redemption" (IMDb 9.3), Streaming: [Prime, Netflix] **</information>**

**<reasoning>** The OTT details tool successfully returned information for "The Shawshank Redemption" **</reasoning>**

**<answer>** The Shawshank Redemption is a critically acclaimed film... It is available on Amazon Prime and Netflix. **</answer>**

**Outcome: Solved.**

Figure 8: Qualitative comparison (Case Study). **Case 2** highlights the robustness of **ToolOmni** against execution failures. While **ToolLlama** gets trapped in a repetitive error loop due to rigid parameter usage, **ToolOmni** demonstrates adaptive planning: after tool failures, it dynamically pivots its strategy—switching from general search to specific verification—to successfully fulfill the user request. *Note: The process is condensed for clarity.*

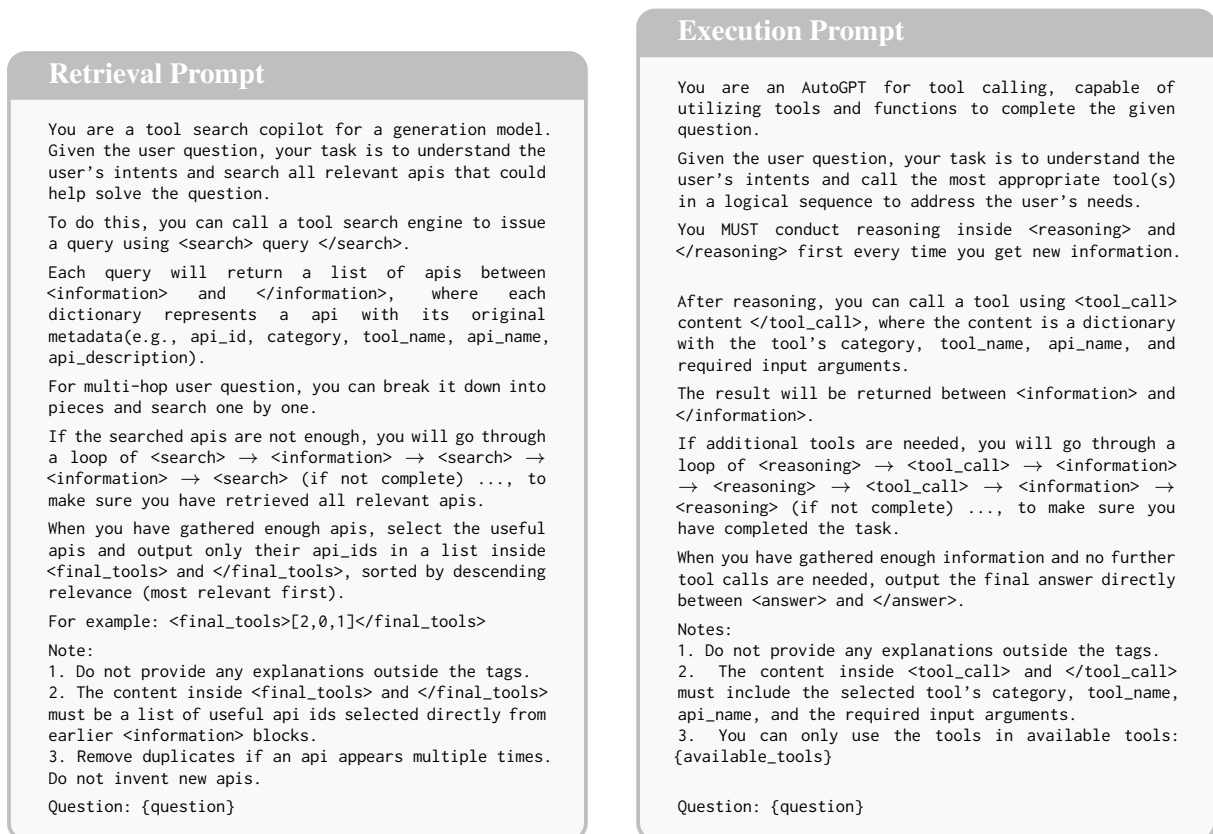


Figure 9: System prompts used for **Retrieval** (Left) and **Execution** (Right) phases. The retrieval prompt guides the agent to proactively search and select tools, while the execution prompt instructs it to perform grounded reasoning and tool invocation.