



CRITICTOOL: Evaluating Self-Critique Capabilities of Large Language Models in Tool-Calling Error Scenarios

Anonymous ACL submission

Abstract

The ability of large language models (LLMs) to utilize external tools has enabled them to tackle an increasingly diverse range of tasks. However, as the tasks become more complex and long-horizon, the intricate tool utilization process may trigger various unexpected errors. Therefore, how to effectively handle such errors, including identifying, diagnosing, and recovering from them, has emerged as a key research direction for advancing tool learning. In this work, we first extensively analyze the types of errors encountered during the function-calling process on several competitive tool evaluation benchmarks. Based on it, we introduce CRITICTOOL, a comprehensive critique evaluation benchmark specialized for tool learning. Building upon a novel evolutionary strategy for dataset construction, CRITICTOOL holds diverse tool-use errors with varying complexities, which better reflects real-world scenarios. We conduct extensive experiments on CRITICTOOL, and validate the generalization and effectiveness of our constructed benchmark strategy. We also provide an in-depth analysis of the tool reflection ability on various LLMs, offering a new perspective on the field of tool learning in LLMs.

1 Introduction

Large Language Models (LLMs) represent a groundbreaking advancement in artificial intelligence, demonstrating remarkable capabilities in various tasks (Zhao et al., 2023; Jiang et al., 2024; Chen et al., 2023; McAleese et al., 2024). The interaction between LLMs and external tools empowers them to address more complex tasks, as these tool-calling systems increasingly adapt to dynamic real-world environments (Chen et al., 2024c).

Driven by practical applications and attractive ability, the evaluation of tool-use capabilities for LLMs remains a topic of ongoing research. Existing works are typically confined to single-tool usage scenarios (Xu et al., 2023; Patil et al., 2023) or

comparing the executions with predefined golden answers (Shen et al., 2023; Ye et al., 2024a,b; Chen et al., 2024b). However, real-world applications often involve complex and multi-step tool-calling tasks, where intricate intermediate trajectories introduce opportunities for errors arising either from LLMs themselves (Yan et al., 2024; Sun et al., 2024) or from external factors (Guo et al., 2024a). Due to the complexity of the external environment, combined with the inherently challenging nature of tool-use tasks, neglecting the process status of tool invocation may result in biased evaluation. Current benchmarks primarily address these challenges by either filtering out erroneous data (Liu et al., 2024) or treating errors as suboptimal nodes to expand the tool answer search space (Qin et al., 2023; Chen et al., 2024a; Abdelaziz et al., 2024). As a result, these approaches fail to provide insights into how LLMs detect and mitigate errors during tool calls, leading to an insufficient evaluation of their tool-use capabilities. Given the diverse sources of errors and the various strategies required to address them, we argue that benchmarks that overlooks LLMs’ error recovery cannot accurately evaluate a model’s actual tool-use performance.

To address these challenges, we introduce CRITICTOOL, the first self-critique evaluation benchmark for tool utilization of LLMs. Distinct from prior result-oriented evaluation methods, we categorize error patterns more finely and evaluate models from multiple perspectives, enabling a deeper exploration of LLMs’ tool-use performance in error-prone scenarios. Specifically, we categorize errors from two main sources: internal model-driven errors and external environment errors. We then diversify our error dataset by ensuring the errors span a wide range of tools and design fine-grained evaluation protocols for two sources of errors. This paradigm enables a granular evaluation of LLMs’ self-critique capabilities across different dimensions: reflect and correct for internal

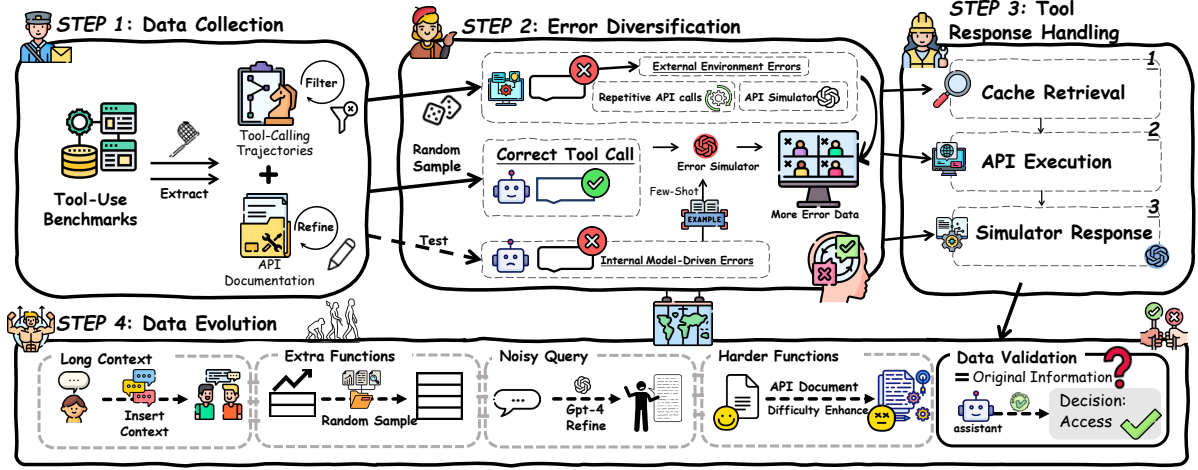


Figure 1: **Overview of CRITICTOOL construction pipeline.** The pipeline begins with collecting and testing tool-use benchmarks to obtain a variety of correct and incorrect tool-calling trajectories. GPT-based simulators and repeated API calls are employed to diversify internal and external error patterns. And responses to internal errors are generated via cache retrieval, API execution, and API simulator. Finally, the error data is evolved using four distinct strategies, followed by verification and manual review.

model-driven errors, and retry with skip or finish for external environment errors.

By conducting extensive experiments on CRITICTOOL, we perform a thorough analysis of the results, providing valuable insights into LLMs’ behavior when encountering different types of errors during tool calls. We observe that different models exhibit varying self-critique behaviors when faced with errors from different sources.

The main contributions of our work are summarized as follows:

- We observe LLMs’ performance in several popular and high-quality tool-use benchmarks and provide a comprehensive analysis of error distributions.
- To the best of our knowledge, we are the first to introduce CRITICTOOL, a tool self-critique evaluation benchmark for LLMs, categorizing errors from different aspects and abilities.
- We propose a novel data evolution strategy to enrich the error dataset by incorporating more complex data scenarios, thus broadening the scope and depth of evaluation for LLMs in real-world applications.
- With extensive experiments, we provide a detailed analysis of the self-critique ability of various LLMs, offering a new perspective in the field of tool learning.

2 CRITICTOOL

In this section, we begin with presenting an in-depth analysis of the key issues in current tool

Table 1: The success rates (%) of GPT-3.5 and Qwen-turbo in recovering from errors across the four datasets.

	Nestful	API-Bank	T-Eval	BFCL
Qwen-turbo	12.64	6.25	35.14	29.47
GPT-3.5	18.10	7.69	51.11	7.14

learning, highlighting the pressing need for tool-specific critique evaluation benchmarks. Building on these observations, we introduce CRITICTOOL, a benchmark designed to systematically explore LLMs’ self-critique capabilities.

2.1 Motivation: LLMs’ Performance on Popular Tool-Use Benchmarks

Tool utilization is a critical yet challenging task in large language model (LLM) applications, requiring sophisticated reasoning and practical adaptation. To identify the current limitations in tool learning, we conduct an in-depth analysis of LLM’s behavioral patterns across various tool-calling benchmarks (Refer to Appendix A for more details). As shown in Tab. 1, our investigation reveals a noteworthy phenomenon: *most LLMs struggle to recover from errors during the tool-calling process, resulting in eventual task failure*. This issue becomes particularly pronounced as tasks grow more complex and long-horizon. Despite the significance of this limitation, existing tool utilization benchmarks rarely directly consider the ability for self-critique, leading to insufficient attention toward improving this capability in tool learning. As highlighted by o1 (OpenAI, 2024), the ability to

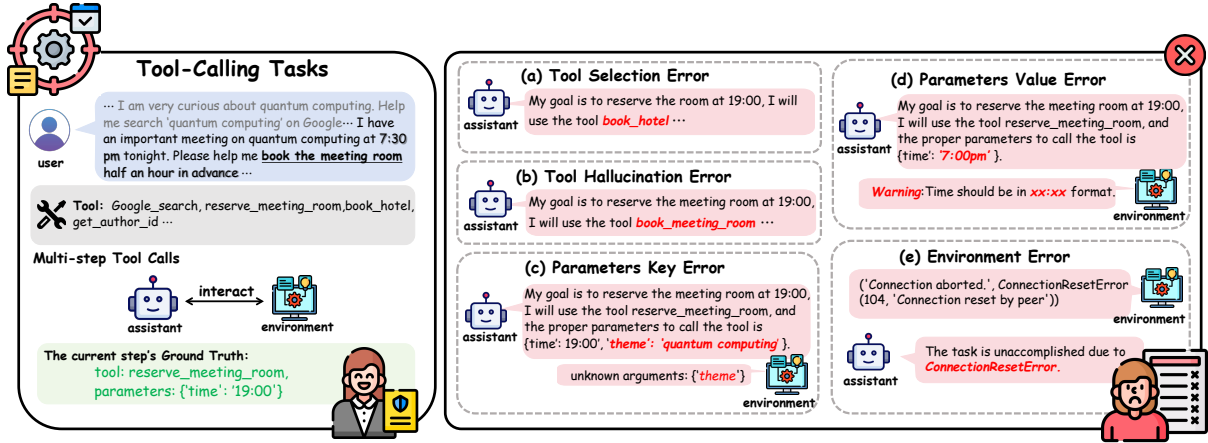


Figure 2: **Examples of Errors in multi-step tool call tasks.** Multi-step tool call errors are categorized into five patterns based on the source and characteristics of the errors: **Tool Selection Errors, Tool Hallucination Errors, Parameters Key Errors, Parameters Value Errors and Environment Errors.**

self-critique is essential for executing long-horizon tasks effectively and serves as a pathway to scalable oversight in LLM reasoning. In this work, we seek to fill this gap by introducing CRITICTOOL, a benchmark designed to systematically evaluate the self-critique capability in tool learning.

2.2 Base Dataset Construction

The construction of the base dataset in CRITICTOOL consists of three main phases: tool-use data collection, error diversification, and tool response handling. The overview of the construction is shown in Fig. 1.

2.2.1 Error Patterns

From our observations of LLMs' tool-use performance in § 2.1, we identify several frequently occurring error patterns when LLMs function as tool-calling assistants, as illustrated in Fig. 2. These errors stem from two primary sources: model capability limitations often give rise to internal model-driven errors related to both tool and parameter handling, while external environment errors will disrupt task completion.

- **Tool Selection Errors:** The assistant selects an existing but unsuitable tool for the given task, often resulting from generating an incorrect goal, or misunderstanding usage of the tool.
- **Tool Hallucination Errors:** The assistant attempts to use a non-existent tool, typically caused by task misinterpretation or failure to recognize available tools.
- **Parameter Key Errors:** The assistant passes incorrect parameter keys, either omitting required ones or including irrelevant keys, usually due to task miscomprehension or forgetting tool require-

ment details.

- **Parameter Value Errors:** The assistant provides incorrect parameter values, usually stemming from failure to comply with the expected input format or overlooking task details.
- **Environment Errors:** Real-world APIs may not always be stable (Guo et al., 2024a). Issues such as connection timeouts or lack of user permissions can disrupt tool interactions, and may cause the assistant to abandon tasks or endlessly retry calls.

2.2.2 Tool-Use Data Collection

To construct CRITICTOOL, our goal is developing a tool-use dataset that spans diverse domains of tools and captures a wide range of errors that LLMs encounter in tool call scenarios. Existing benchmarks have already collected realistic APIs and generated well-designed tool-use tasks with excellent diversity and appropriate complexity, making them ideal sources of tool-use data. We use the datasets from high-quality tool-use benchmarks, including BFCL v3 (Yan et al., 2024) and T-Eval (Chen et al., 2024b), which provide access to 203 real-world APIs across 23 tools and a variety of multi-step tool-use tasks that require complex agent-tool interactions, perfectly aligning with our goals.

We have curated error-containing data while observing LLMs' behavioral patterns across these benchmarks in § 2.1, but it is far from sufficient. To facilitate more controlled error data generation, we first collect the ground truth tool-calling trajectories including tool call actions and the corresponding tool responses across various tasks in these datasets. Any data containing errors, such as incorrect annotations or failed tool calls, is carefully manually filtered to ensure the quality and reliability of our

dataset. Next, we extract API documentation and refine any ambiguous or inadequate descriptions to ensure clarity and precision, minimizing potential misunderstandings. To further enhance consistency, we standardize all tool-calling trajectories and API descriptions, which aligns formats across different benchmarks, creating a coherent framework that facilitates consistent prompts and reliable tool-use interactions throughout our evaluation. Examples are provided in the Appendix C.1.

2.2.3 Error Diversification

We have identified five patterns of errors from two sources in § 2.2.1. To ensure the comprehensive coverage of potential scenarios, we systematically diversify these errors, significantly expanding our error repository.

- **Internal Model-Driven Errors:** The internal model-driven error data collected from previous observation has two limitations that (1) it comes from a small subset of tools and tasks, and (2) the tests primarily involve advanced LLMs, which restricts the coverage of errors that less capable models might produce. Moreover, our observation reveals that LLMs tend to exhibit similar behaviors within a specific error pattern, despite interacting with different tools. This similarity allows us to expand the diversity of errors in the calling of all tools. We prompt GPT-4o as an error simulator, simulating error-prone behaviors of tool-calling assistants. Using examples of error patterns collected from observation as few-shot demonstrations (Brown et al., 2020), error simulator is tasked with generating diverse instances of errors across a wider range of tools and tasks.

- **External Environment Errors:** During data collection, we capture numerous instances of tool responses containing external environment errors and match them with their corresponding tools. However, not all tools in the benchmark datasets include such error examples. To fill this gap, we implement repeated API calls and API simulation strategies. We perform repeated calls to the accessible APIs to collect the error responses, which occur due to the inherent instability of the external environment. For the inaccessible APIs, we also employ GPT-4o as an API simulator to collect some environment error responses.

2.2.4 Tool Response Handling

The responses LLMs receive from the environments during tool calls are crucial for them to

self-criticize, making it essential to obtain tool responses corresponding to internal model-driven errors. However, due to permission restrictions, not all collected APIs are executable. Inspired by StableToolBench (Guo et al., 2024a), we adopt a systematic approach for tool response collection based on the availability status of each API.

- **Cache Retrieval:** We first search the cache to check whether the tool and parameters used in the current call have previously been cached. If a match is found, the cached response is used as the environment’s response for the current tool call.

- **API Execution:** If there is no match in the cache, we then verify the accessibility of API. The tool call is executed and the actual API response is used if the API is available.

- **Simulator Response:** When neither cache nor API is available, we employ GPT-4o as an API simulator to ensure that the tool-calling assistant still receives feedback for its current action.

2.3 Data Evolution

Real-world tool calls typically encompass complex contexts, sophisticated tools, and ambiguous user queries (Wang et al., 2024b). To achieve a more realistic evaluation of LLM performance in tool call tasks, we propose a strategy termed Scalable and Robust Mixed Self-Evolution (SRM) to facilitate the self-evolution of data within the origin benchmark. Specifically, we focus on two critical factors of tool-use tasks: scale and robustness. Based on these factors, we develop four distinct evolutionary sub-strategies on these perspectives that closely align LLM tool-use tasks with real-world scenarios while preserving the ground truth.

- **Long Context:** We introduce extended conversations that range from 1k to 3k tokens from LongBench (Bai et al., 2023) as the context and randomly insert them prior to the user’s tool-use query.

- **Extra Tools:** Most existing benchmarks merely supply the tools required for specific test tasks, which contrasts sharply with the vast number of APIs involved in real applications. Thus, we propose the Extra Tools evolution strategy, which randomly incorporates additional tools into API lists.

- **Noisy Query:** Real user queries are often verbose, vague, include unnecessary information, and are prone to typographical errors, which challenge LLMs’ ability to interpret intent. We employ GPT-4o to simulate human language habits, particular focusing on addressing irrelevant information, cumbersome expressions, and typographical issues.

• **Harder Tools:** DRAFT (Qu et al., 2024) and BFCL v2 (Yan et al., 2024) illustrate the substantial impact that API documentation has on LLM tool calls. Therefore, we deliberately degrade the API document by prompting GPT-4o, thereby making the idealized APIs documentation more realistic.

We combine the four evolutionary sub-strategies to increase the difficulty of LLM tool-use tasks, involving three key components: context, queries, and the API list, enabling the exploration of scalability and robustness in self-critique.

After the SRM process, we verify the data to ensure that the ground truth remains unchanged. Due to the evolution process, it is difficult to determine whether inappropriate self-critique behavior in following evaluation arises from the model’s inherent limitations or biases introduced by the evolutionary strategies. Moreover, re-annotating the evolved data when ground truth is available would neither be cost-effective nor environmentally sustainable. To address this, we devise a novel data verification approach, termed equivalence verification. We use GPT-4o to check whether the modifications or additions made during the evolution process significantly impact the tool-use tasks. The specific implementation details are provided in the Appendix C.2. Finally, human experts are employed for double-check.

2.4 Fine-Grained Evaluation

CRITICTOOL comprehensively evaluates the self-critique capabilities of LLMs by breaking them down into multiple dimensions, across different error patterns encountered during tool interaction.

2.4.1 Self-Critique Task Decomposition

In the CRITICTOOL, each tool-use task is defined as a tuple (Q, T) , where Q is the task query, and T represents the list of APIs available for the tool-calling assistant. We define the trajectory \mathcal{T} as a sequence of tool-response pairs $\{(a_i, r_i)\}$, capturing the interaction between the assistant’s action a and the corresponding tool response r in the i -th step. The action a is regarded as either $(goal, tool, args)$ or $(tool, args)$ depending on whether the chain of thought strategy is applied.

The complex interactions between the assistant and the environment can lead to potential errors at any step, underscoring the importance of evaluating LLMs’ self-critique capabilities at the step level (Ye et al., 2024b). Consequently, the test data consists of the first k steps of the tool-calling trajec-

tory for each task, where k is randomly chosen, and any errors may be introduced at step k . For evaluation, we define the solution path $S = (c, \hat{a})$, where c represents the critique of the error when the tool call action a_k contains an error, and $S = (\hat{a})$ or a sequence of actions $S = \{\hat{a}_1, \hat{a}_2, \dots\}$ otherwise.

In tasks evaluating self-critique abilities for internal model-driven errors, CRITICTOOL employs both error-free and error-injected data to ensure fairness and robustness. We evaluate the $(k + 1)$ -th step and deconstruct the self-critique process into two dimensions. The tool-calling assistant should recognize whether an error occurred during the preceding tool call first and identify its specific category. This process of identifying and analyzing errors is defined as **reflect**, a central step in the model’s self-critique. Based on the result of the reflection, the model needs to take corrective action to recover from the error. We define this process as **correct**, highlighting the model’s ability to improve and adapt its behavior effectively.

For tasks involving external environment errors, the assistant is expected to properly handle the response from the environment that contains the error signal in the subsequent steps. We encourage the assistant to **retry** the failed tool calls a limited number of times to avoid incidental error caused by environmental instability. If the issue persists despite multiple retries, the assistant should **skip** the problematic step and address any remaining feasible subtasks or **finish** the tool-calling process and inform the user that further guidance is required.

2.4.2 Evaluation Metrics

CRITICTOOL employs fine-grained evaluation metrics to assess each dimension of self-critique behavior of LLMs across different error scenarios. The details can be found in Appendix C.3.

• **REFLECT:** The reflect evaluator asks the assistant to determine whether to produce a critique c^{pred} , based on the correctness of tool call action a_k . Then, c^{pred} is compared with the golden answer c^{gt} if an error exists in a_k .

• **CORRECT:** The correct evaluator asks the assistant to generate a corrected action \hat{a}^{pred} for a detected error in tool call action a_k , and compares \hat{a}^{pred} with the golden answer \hat{a}^{gt} .

• **RETRY:** The assistant is asked to generate a repeated tool call \hat{a}_1^{pred} if any error signal is found in r_k . The evaluator compares \hat{a}_1^{pred} with the golden answer \hat{a}_1^{gt} , which corresponds to the action a_k .

• **SKIP:** If the error from the environment can-

Table 2: **Main Results of CRITICTOOL.** **Bold** indicates the best performance across all models, while underline denotes the best performance within the same group and scale of models.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect	Category	Correct	Args	Retry	Break	Skip/Finish	Args	
	Detect		Tool				Tool		
<i>Close-Source Large Language Models</i>									
Claude3.5	<u>82.21</u>	56.28	84.52	77.70	<u>38.80</u>	56.37	22.29	26.39	55.88
GPT-3.5	75.43	63.37	70.18	55.21	10.85	89.27	51.65	43.39	60.88
GPT-4o	79.53	<u>71.18</u>	85.52	80.13	18.51	<u>96.46</u>	<u>52.83</u>	43.62	69.78
<i>Open-Source Large Language Models</i>									
LLaMA3-8B	49.01	31.63	67.36	61.39	36.78	73.53	31.93	<u>30.01</u>	49.54
LLaMA3.1-8B	<u>84.72</u>	68.32	78.79	69.93	50.94	78.18	26.77	23.63	59.45
Qwen2.5-7B	83.64	43.68	77.26	69.17	29.20	88.23	<u>40.64</u>	22.62	58.88
GLM4-9B-chat	57.51	25.34	60.51	50.22	19.22	<u>90.45</u>	36.56	23.02	48.36
Minstral-8B	48.62	26.20	68.89	59.26	49.76	50.47	17.45	20.81	42.50
LLaMA3-70B	56.65	29.94	69.41	62.91	33.02	73.52	28.02	27.81	49.56
LLaMA3.1-70B	81.65	61.62	82.77	66.99	65.40	92.63	54.16	27.32	66.18
Qwen2.5-72B	86.83	55.35	<u>83.36</u>	<u>76.85</u>	40.68	97.05	55.54	<u>32.91</u>	<u>68.11</u>
<i>Tool-Use-Finetuned Large Language Models</i>									
ToolLLaMA2-7B	0.95	0.00	3.99	0.75	0.84	0.90	0.74	0.00	1.10
ToolACE-8B	12.61	<u>0.88</u>	13.01	11.78	1.30	17.92	8.73	13.59	10.17
AgentLM-7B	<u>22.93</u>	0.00	<u>46.91</u>	<u>36.47</u>	<u>11.81</u>	<u>81.38</u>	<u>18.80</u>	<u>17.64</u>	<u>33.06</u>

not be resolved within the retry limit, the assistant should skip and proceed with the next feasible subtask. The skip action \hat{a}_n^{pred} is compared to the golden answer \hat{a}_2^{gt} , which indicates the ground truth action for the next subtask.

- **FINISH:** The evaluator checks whether the assistant terminates the tool call and waits for further instructions from the user after several unsuccessful attempts to resolve the environmental error.

- **OVERALL:** We calculate the overall score by weighing the self-critique dimensions based on their importance in completing a tool-calling task. The weight assigned to reflect is 0.2, to correct is 0.3, to retry is 0.05, and to skip/finish is 0.45.

3 Experiment

3.1 Experiment Setup

We conduct evaluations on CRITICTOOL using a diverse set of 14 LLMs, to establish a comprehensive self-critique benchmark for assessing the capabilities of current large language models. For closed-source LLMs, we select three prominent models: Claude3.5 (Anthropic, 2024) developed by Anthropic, alongside GPT-3.5 (OpenAI, 2022) and GPT-4o (Hurst et al., 2024) provided by OpenAI.¹ For open-source LLMs, we evaluate numerous models including LLaMA3, LLaMA3.1 (AI@Meta, 2024), Qwen2.5 (Team,

¹The version for GPT-4o is gpt-4o-2024-08-06, for GPT-3.5 is gpt-3.5-turbo-16k, and for Claude3.5 is claude-3-5-sonnet-20241022.

2024a,b), GLM4 (GLM et al., 2024), Minstral (AI, 2024). For tool-use-finetuned LLMs, we evaluate ToolLLaMA2 (Qin et al., 2023), ToolACE (Liu et al., 2024) and AgentLM (Zeng et al., 2023).

3.2 Benchmarking Results on CRITICTOOL

The detailed experimental results are shown in Tab. 2. Experiments using the chain of thought strategy (Wei et al., 2022) are also conducted, leading to improvements in LLMs’ self-critique performance, with the results provided in the Appendix D.3. We analyze the benchmarking results by exploring the following four questions.

Q1: Which Model is Better at Tool Self-Critique?

GPT-4o leads in self-critique performance for tool-use error scenarios, achieving an impressive overall score of 69.78. Close behind, large-scale open-source models LLaMA3.1-70B and Qwen2.5-72B, deliver comparable scores, showcasing strong self-critique capabilities.

For internal model-driven errors, the closed-source models GPT-4o and Claude3.5 deliver comparable top performances, though Claude3.5 slightly underperforms in error categorization. In contrast, open-source models exhibit substantial variability in self-critique performance. While most open-source models significantly lag behind the closed-source models, highlighting a clear gap in their capabilities, LLaMA3.1 and Qwen2.5 stand out as notable exceptions. Their performance

not only approaches but occasionally surpasses that of closed-source models. However, tool-use-finetuned models show disappointing results in handling internal errors. Except for AgentLM-8B, the other models exhibit almost no instruction-following or self-critique capabilities, which can be attributed to the damage to their generalization ability caused by fine-tuning on specific data.

For external environment errors, most models can recognize errors and avoid endless repetition, though Claude3.5 and Minstral-8B shows weaker performance in this regard, and some tool-use-finetuned models entirely lack this ability. When it comes to handling errors by either proceeding with subsequent tasks or finish tool call action, GPT-4o outperforms other models, with some large-scale open-source models achieving comparably strong performance.

Q2: How do Models Perform in Self-Critique across Different Internal Error Patterns?

As shown in Fig. 3, we analyze self-critique performance on internal error patterns by focusing on GPT-4o and LLaMA3.1-8B, the strongest close-source and small-scale open-source model. The results of more models can be found in Appendix D.

Tool selection errors, often manifesting as silent errors without clear external signals (Sun et al., 2024), are the most challenging error for model to detect, resulting in low reflect accuracy and poor correction performance across models. In contrast, tool hallucination errors are easier to detect due to their more evident inconsistencies. GPT-4o demonstrates a clear advantage in reflecting on such errors compared to LLaMA3.1-8B and other open-source models. Both models exhibit high reflect accuracy for parameters key and value errors, with parameters key errors being relatively easier to correct.

Overall, models are better at reflecting on errors with clear external signals. Furthermore, correcting tool-related errors is inherently more complex, as it involves ensuring both correct tool selection and accurate parameters passing. Consequently, parameters-related errors, which require only adjustments to the passed parameters, are corrected with higher accuracy.

Q3: How does Data Evolution Effects?

As illustrated in Fig. 4, the SRM strategy leads to a decline in the scores of all LLMs. GPT-4o retains its SOTA results, while LLaMA3.1-8B and Qwen2.5-7B also demonstrate impressive capabilities. In contrast, LLaMA3-70B experiences significant performance degradation, falling below the

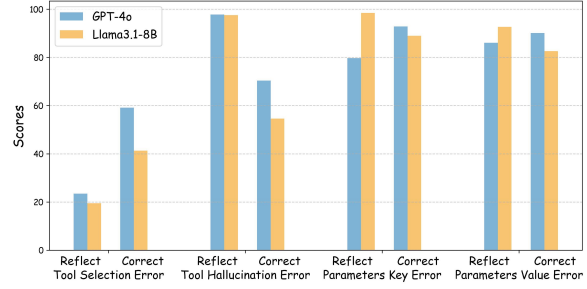


Figure 3: Comparison of scores between GPT-4o and Llama3.1-8B.

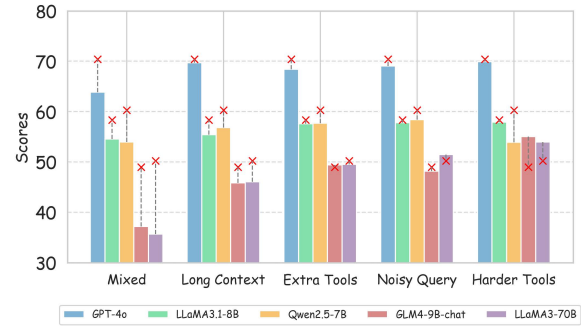


Figure 4: Comparison of the performance of five models across various evolution strategies. The red cross indicates the score corresponding to the base dataset.

performance of most small scale models. This is consistent with CriticBench (Lin et al., 2024) experimental observation. We attribute this to the unstable generalizability of the offline data, a limitation that becomes increasingly pronounced as the number of model parameters grows. We independently test the four sub-strategies to investigate their impact on models’ self-critic performance. The differences in the negative impact of the four evolutionary sub-strategies are not significant, and the negative impact on the model decreases in the following order: Long Context, Extra Tools, Noisy Query, and Harder Tools. Long Context and Extra Tools increase the difficulty of retrieval and challenge the model’s ability to follow instructions. Noisy Query and Harder Tools do not introduce excessive additional information, but diminish the LLM model’s ability to understand different tools and verbose user queries. However, as the API documents become more verbose and longer, some models demonstrate improved comprehension of the APIs, leading to slight performance enhancements, such as GLM4-9B-chat.

Overall, for the model, the three key components—the context, query, and tool list—are not merely superimposed. The interplay between scalable and robust levels results in a compounding

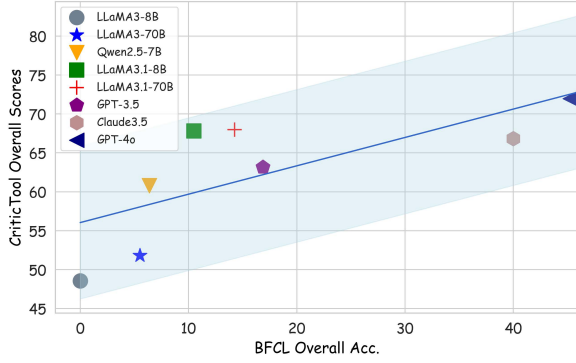


Figure 5: Comparison between BFCL Overall Accuracy and CRITICTOOL Overall Scores across several models. LLMs show similar trends in tool-use and self-critique capabilities.

effect, causing the model’s performance to degrade more rapidly under the hybrid strategy compared to individual strategies. The detailed results can be found in Appendix C.2.3.

Q4: What is the Relationship Between Tool-Use and Self-Critique Capabilities?

We compare the fine-grained evaluations on CRITICTOOL with the results of the benchmark designed to explore tool-use capabilities, investigating the relationship between models’ self-critique capabilities in tool-calling tasks and their tool-use capabilities. We analyze the Overall Accuracy metric from the BFCL v3 (Yan et al., 2024), which includes multi-step tool-calling scenarios, to examine the relationship between the tool-use performances of selected models and their Overall performance on CRITICTOOL. As results shown in Fig. 5, we observe a general alignment between the trends in models’ tool-use and self-critique capabilities. This observation not only indicates a strong connection between models’ ability to accurately use tools and their self-critique capabilities, suggesting that strengthening self-critique mechanisms could provide a promising avenue for enhancing overall tool-use performance, but also validates the rationale behind our benchmark.

4 Related Work

Tool Learning with LLM There are currently two primary technical approaches for enhancing the tool invocation capability of LLMs (Shen et al., 2023; Yuan et al., 2024). The first approach focuses on constructing high-quality tool call data and improving the model’s tool invocation capabilities through fine-tuning (Kong et al., 2024; Chen et al., 2024a; Patil et al., 2023). The second ap-

proach involves leveraging contextual tool invocation demonstrations to augment the model’s ability to invoke tools through in-context learning (Wang et al., 2024a).

The evaluation of tool invocation capabilities across different models is also an urgent issue. Common evaluation frameworks involve comparing model predictions to ground truth (Yan et al., 2024; Guo et al., 2024b), while ToolBench (Qin et al., 2023) contrasts model predictions with those generated by advanced LLMs, such as GPT-4. Although some studies (Yan et al., 2024; Yao et al., 2024; Sun et al., 2024) have identified common errors in tool invocations, they unfortunately lack in-depth analysis and the design of targeted evaluation frameworks. In contrast to the aforementioned benchmarks, CRITICTOOL is the first to analyze various errors and evaluate the self-critic ability in tool invocation as far as we know.

Self-Critique of LLMs Learning from incorrect attempts can help prevent similar errors, thereby enabling deeper insights into the data and facilitating self-learning (Ke et al., 2024; Shinn et al., 2023; An et al., 2023; Ying et al., 2024; Zhang et al., 2024; Tian et al., 2024). CriticEval (Lan et al., 2024) evaluate the self-critique ability of LLMs on nine key tasks, including math and code, across four critical dimensions. For tool calls, the self-critic strategy is particularly well-suited for this complex task, which integrates various important capabilities on massive and constantly updated tools (Gou et al., 2023). However, to the best of our knowledge, no prior work has specifically explored the evaluation of self-critique in tool invocations. Recognizing the unique characteristics of tool calls compared to other tasks, CRITICTOOL adopts a targeted and fine-grained evaluation framework.

5 Conclusion

In this paper, we propose CRITICTOOL, the first benchmark for tool self-critique in LLM tool evaluation as far as we know. CRITICTOOL explicitly distinguishes between internal model errors and external environment errors, classifies evaluation methods, and employs data evolution strategies to uncover the true capabilities of the models under evaluation. This evaluation offers a comprehensive analysis and identifies the primary bottlenecks in current LLMs’ tool learning, providing valuable insights for the future development of tool agents.

References

- Ibrahim Abdelaziz, Kinjal Basu, Mayank Agarwal, Sadhana Kumaravel, Matthew Stallone, Rameswar Panda, Yara Rizk, GP Bhargav, Maxwell Crouse, Chulaka Gunasekara, et al. 2024. Granite-function calling model: Introducing function calling abilities via multi-task learning of granular tasks. In *EMNLP*, pages 1131–1139.
- Mistral AI. 2024. *Un ministral, des ministraux*.
- AI@Meta. 2024. *Llama 3 model card*.
- Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng, Jian-Guang Lou, and Weizhu Chen. 2023. Learning from mistakes makes llm better reasoner. *arXiv preprint arXiv:2310.20689*.
- Anthropic. 2024. *Claude 3.5 sonnet*.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.
- Kinjal Basu, Ibrahim Abdelaziz, Kelsey Bradford, Maxwell Crouse, Kiran Kate, Sadhana Kumaravel, Saurabh Goyal, Asim Munawar, Yara Rizk, Xin Wang, et al. 2024. Nestful: A benchmark for evaluating llms on nested sequences of api calls. *arXiv preprint arXiv:2409.03797*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *NeurIPS*.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*.
- Sijia Chen, Yibo Wang, Yi-Feng Wu, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Lijun Zhang. 2024a. Advancing tool-augmented large language models: Integrating insights from errors in inference trees. *arXiv preprint arXiv:2406.07115*.
- Zehui Chen, Weihua Du, Wenwei Zhang, Kuikun Liu, Jiangning Liu, Miao Zheng, Jingming Zhuo, Songyang Zhang, Dahua Lin, Kai Chen, and Feng Zhao. 2024b. T-eval: Evaluating the tool utilization capability of large language models step by step. In *ACL*, pages 9510–9529.
- Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. 2024c. Agent-FLAN: Designing data and methods of effective agent tuning for large language models. In *ACL*, pages 9354–9366.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2023. Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*.
- Zhicheng Guo, Sijie Cheng, Hao Wang, Shihao Liang, Yujia Qin, Peng Li, Zhiyuan Liu, Maosong Sun, and Yang Liu. 2024a. StableToolBench: Towards stable large-scale benchmarking on tool learning of large language models. In *ACL*, pages 11143–11156.
- Zishan Guo, Yufei Huang, and Deyi Xiong. 2024b. CToolEval: A Chinese benchmark for LLM-powered agent evaluation in real-world API interactions. In *ACL*, pages 15711–15724.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. 2024. A survey on large language models for code generation. *arXiv preprint arXiv:2406.00515*.
- Pei Ke, Bosi Wen, Andrew Feng, Xiao Liu, Xuanyu Lei, Jiale Cheng, Shengyuan Wang, Aohan Zeng, Yuxiao Dong, Hongning Wang, Jie Tang, and Minlie Huang. 2024. CritiqueLLM: Towards an informative critique generation model for evaluation of large language model generation. In *ACL*, pages 13034–13054.
- Yilun Kong, Jingqing Ruan, YiHong Chen, Bin Zhang, Tianpeng Bao, Shi Shiwei, du Guo Qing, Xiaoru Hu, Hangyu Mao, Ziyue Li, Xingyu Zeng, Rui Zhao, and Xueqian Wang. 2024. TPTU-v2: Boosting task planning and tool usage of large language model-based agents in real-world industry systems. In *EMNLP*, pages 371–385.
- Tian Lan, Wenwei Zhang, Chen Xu, Heyan Huang, Dahua Lin, Kai Chen, and Xian-ling Mao. 2024. Criticbench: Evaluating large language models as critic. *arXiv preprint arXiv:2402.13764*.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. API-bank: A comprehensive benchmark for tool-augmented LLMs. In *EMNLP*, pages 3102–3116.
- Zicheng Lin, Zhibin Gou, Tian Liang, Ruilin Luo, Haowei Liu, and Yujia Yang. 2024. CriticBench: Benchmarking LLMs for critique-correct reasoning. In *ACL*, pages 1552–1587.

Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, et al. 2024. Toolace: Winning the points of llm function calling. <i>arXiv preprint arXiv:2409.00920</i> .	788
Nat McAleese, Rai Michael Pokorny, Juan Felipe Ceron Uribe, Evgenia Nitishinskaya, Maja Trebacz, and Jan Leike. 2024. Llm critics help catch llm bugs. <i>arXiv preprint arXiv:2407.00215</i> .	793
OpenAI. 2022. Introducing chatgpt .	794
OpenAI. 2024. Introducing openai o1 .	795
Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. Gorilla: Large language model connected with massive apis. <i>arXiv preprint arXiv:2305.15334</i> .	796
Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. <i>arXiv preprint arXiv:2307.16789</i> .	797
Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2024. From exploration to mastery: Enabling llms to master tools via self-driven interactions. <i>arXiv preprint arXiv:2410.08197</i> .	798
Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .	799
Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. 2022. Fine-tuned language models are continual learners. In <i>EMNLP</i> , pages 6107–6122.	800
Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. 2023. Taskbench: Benchmarking large language models for task automation. <i>arXiv preprint arXiv:2311.18760</i> .	801
Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In <i>NeurIPS</i> .	802
Jimin Sun, So Yeon Min, Yingshan Chang, and Yonatan Bisk. 2024. Tools fail: Detecting silent errors in faulty tools. In <i>EMNLP</i> , pages 14272–14289.	803
Qwen Team. 2024a. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> .	804
Qwen Team. 2024b. Qwen2.5: A party of foundation models .	805
Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. 2024. Toward self-improvement of llms via imagination, searching, and criticizing. <i>arXiv preprint arXiv:2404.12253</i> .	806
Boshi Wang, Hao Fang, Jason Eisner, Benjamin Van Durme, and Yu Su. 2024a. LLMs in the imagerie: Tool learning through simulated trial and error. In <i>ACL</i> , pages 10583–10604.	807
Siyuan Wang, Zhuohan Long, Zhihao Fan, Zhongyu Wei, and Xuanjing Huang. 2024b. Benchmark self-evolving: A multi-agent framework for dynamic llm evaluation. <i>arXiv preprint arXiv:2402.11443</i> .	808
Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>NeurIPS</i> .	809
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. <i>arXiv preprint arXiv:2305.16504</i> .	810
Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. 2024. Berkeley function calling leaderboard.	811
Jihan Yao, Wenxuan Ding, Shangbin Feng, Lucy Lu Wang, and Yulia Tsvetkov. 2024. Varying shades of wrong: Aligning llms with wrong answers only. <i>arXiv preprint arXiv:2410.11055</i> .	812
Junjie Ye, Guanyu Li, Songyang Gao, Caishuang Huang, Yilong Wu, Sixian Li, Xiaoran Fan, Shihan Dou, Qi Zhang, Tao Gui, et al. 2024a. Tooleyes: Fine-grained evaluation for tool learning capabilities of large language models in real-world scenarios. <i>arXiv preprint arXiv:2401.00741</i> .	813
Junjie Ye, Yilong Wu, Songyang Gao, Caishuang Huang, Sixian Li, Guanyu Li, Xiaoran Fan, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024b. RoTBench: A multi-level benchmark for evaluating the robustness of large language models in tool learning. In <i>EMNLP</i> , pages 313–333.	814
Jiahao Ying, Mingbao Lin, Yixin Cao, Wei Tang, Bo Wang, Qianru Sun, Xuanjing Huang, and Shuicheng Yan. 2024. LLMs-as-instructors: Learning from errors toward automating model improvement. In <i>EMNLP</i> , pages 11185–11208.	815
Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Yongliang Shen, Ren Kan, Dongsheng Li, and Deqing Yang. 2024. Easytool: Enhancing llm-based agents with concise tool instruction. <i>arXiv preprint arXiv:2401.06201</i> .	816
Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2023. Agenttuning: Enabling generalized agent abilities for llms. <i>arXiv preprint arXiv:2310.12823</i> .	817
Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming Lu. 2024. Self-contrast: Better reflection through inconsistent solving perspectives. In <i>ACL</i> , pages 3602–3622.	818

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen
Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang,
Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu,
Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023.
A survey of large language models. *arXiv preprint*
arXiv:2303.18223.

A Observation: Insight into LLMs' Tool-Use Performance

In § 2.1, we test BFCL v3 (Yan et al., 2024), T-Eval (Chen et al., 2024b), API-Bank (Li et al., 2023), and NESTFUL (Basu et al., 2024) to conduct an in-depth analysis of LLMs' behavioral patterns. The details of these benchmarks are provided below.

BFCL V3 is a comprehensive benchmark for evaluating LLMs' performance in multi-step and multi-turn tool calling. The benchmark includes 200 basic tool-use trajectories, along with an additional 800 trajectories that introduce various complexities built upon these basic data.

T-Eval provides 553 tool-use trajectories, breaking down tasks into sub-processes including instruction following, planning, reasoning, retrieval, understanding, and review.

API-bank has 314 tool-use trajectories to evaluate LLMs' capabilities in planning, retrieving, and calling APIs.

NESTFUL is designed to better evaluate LLMs on nested sequences of tool calls. It compiles 85 executable tool-use traces and 215 non-executable traces from the different datasets, as well as synthetic data generated by LLMs.

We first observe that the prompts and tool-call formats used in these benchmarks varied, which could lead to discrepancies in how LLMs follow instructions. To address this, we standardize the test data into a consistent format, as Fig. 11, ensuring LLMs execute tasks sequentially and consistently across benchmarks. Then, we randomly select a subset of the test data from these benchmarks and summarize the frequently occurring error patterns in the test results. The distribution of error patterns is shown in Tab. 3. In the experiment, we observe LLMs' performance in the presence of errors, and gain insight into their different behavior across different errors, as shown in Fig. 12 and 13. When LLMs continue executing tool-use tasks after making mistakes, we find that some of them could recognize and correct their mistakes, while most perform poorly. In cases where tool responses contain errors due to instability, many LLMs become trapped in repetitive retry loops, with few capable of recognizing the issue and breaking free by either skipping the current step or terminating the task.

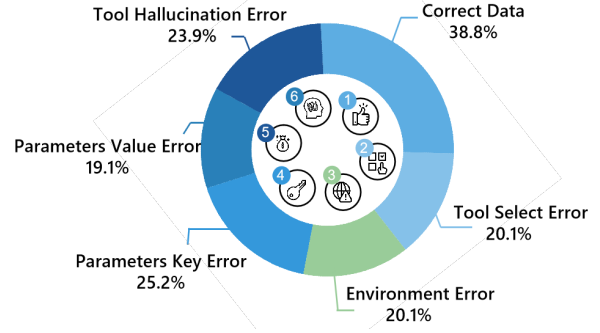


Figure 6: Error distribution for Base data in CRITICTOOL.

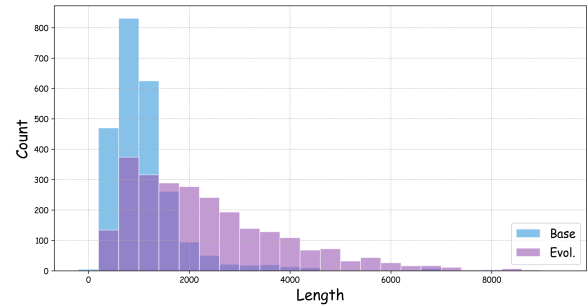


Figure 7: Length distribution for Base and Evolution data in CRITICTOOL, measured by the number of tokens.

B CRITICTOOL Benchmark Details

B.1 Dataset Summary

The base dataset of CRITICTOOL originates from 733 high-quality tool-call trajectories, consisting of 1490 test cases in total, which contains 1316 internal model-driven error test cases and 174 external environment error test cases. On this basis, we retain the error distribution on the base data and randomly select to construct CRITICTOOL evolution dataset (be simplified to Evol.), generating 1000 internal and 250 external new test cases. We visualize the error distribution and length distribution for the base and evol datasets.

Fig. 6 illustrates the error distribution of CRITICTOOL, which comprehensively covers the behavior patterns of LLMs observed across mainstream benchmarks.

Fig. 7 shows that each set of the base benchmark has 1291 tokens on average, while each Evol. contains 2387 tokens on average, validating the generalization and discrimination for tool utilization self-critic evaluation.

Table 3: Error distribution among LLMs in tool-use benchmarks.

Benchmark	Model	Total	Tool Sel. Errors	Tool Halluc. Errors	Param. Key Errors	Param. Value Errors
BFCL V3	GPT-3.5	202	85	0	0	13
	Qwen-turbo	184	82	1	0	13
T-Eval	GPT-3.5	466	38	13	10	29
	Qwen-turbo	452	36	3	4	36
API-bank	GPT-3.5	275	6	1	1	18
	Qwen-turbo	259	2	1	0	13
NESTFUL	GPT-3.5	215	13	22	20	22
	Qwen-turbo	215	9	1	27	29

C Implementation Details

C.1 Data Collection

We collect 733 ground truth tool-calling trajectories from high-quality tool-use benchmarks, BFCL (Yan et al., 2024) and T-Eval (Chen et al., 2024b). To facilitate following controlled error data generation, we manually filter out 485 trajectories that contain no errors and refine the API documentation to ensure that all API descriptions are clear and accurate. To bridge the gap between different instruction formats, we standardize both the trajectories and API documentation, as illustrated in Fig. 14 and 15. This standardization ensures compatibility and reduces variability in the data, enabling a more consistent evaluation of LLMs’ performance in self-critique capabilities.

C.2 Prompts Demonstration

Refer to the corresponding prompt block for a detailed demonstration.

C.2.1 Error Data Diversification

We prompt GPT-4o as error simulator, and the corresponding prompt is presented in Fig. 16.

C.2.2 Tool Responses Generation

We prompt GPT-4o as API simulator, and the corresponding prompt is presented in Fig. 17.

C.2.3 Data Evolution

The framework of the data evolution has been shown in Fig. 9. And Tab. 4, presents a simplified example of our Scalable and Robust Mixed Self-Evolution(SRM) evolution strategy.

Long Context: Recent work (Liu et al., 2024) has demonstrated the importance of context for

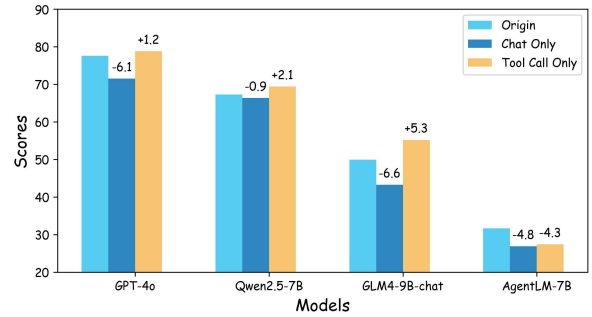


Figure 8: Comparison of scores of different models in Chat Only and Tool Call Only as the context.

error recovery in tool invocations. To this end, we replace the data extracted from the Long Context in § 2.3 with the previously filtered data (including both correct and incorrect samples) and conduct a comparative experiment after manually ensuring no overlap between tasks and tools. As shown in Fig.8, the scores of most models improved to some extent on Only Tool Call. We argue that the tool call context provides a few-shot format for recovery, functioning similarly to an experience replay strategy (Scialom et al., 2022). Therefore, to eliminate unnecessary influence, we rely solely on pure dialogue as the source for Long Context Evolution.

Noisy Query: We prompt GPT-4o to downgrade the API document, and the corresponding prompt is presented in Fig. 18.

Harder Tools: We prompt GPT-4o to downgrade the API documentation, and the corresponding prompt is presented in Fig. 19.

Data Verification: We prompt GPT-4o to verify the evolution data, and the corresponding prompt is presented in Fig. 20, 21, 22, 23.

Table 4: A simplified example of our data evolution strategy.

Original Tool Call Trajectory			
Context: None.			
Tool List: 'name': 'ReserveMeeting_get_room_status', 'description': 'a Tool that get the room booking status'			
User Query: Could you check if there are any available meeting rooms between 14:00 and 16:00?			
Ground Truth: ReserveMeeting_get_empty_room_time(rooms: '[103]')			
Perspective	Sub-strategy	Changed Items	Examples
Scalable	Long Context	Context	Insert Context: We are convening a meeting to review and strategize on our ongoing project. This gathering is crucial for aligning our efforts and ensuring collective success. Your presence is vital as we chart the project's trajectory.
	Extra Tools	Tool List	Extra Tools: Email.show, Email.send, Email.read, ArxivSearch.get_arxiv_information, BINGMap.search_nearby...
Robust	Noisy Query	User Query	Refine Query: Whether it would be possible for you to take a moment to verify if there are any meeting rooms that happen to be unoccupied or not in use between the hours of 2:00 in the afternoon and 4:00 in the afternoon.
	Harder Tools	Tool List	Refine API Document: get <i>rom(room)</i> status

C.3 Detailed Evaluation Metrics

In the CRITICTOOL, self-critique capabilities are divided into multiple dimensions based on errors from different sources: Reflect, Correct, Retry, and Skip/Finish. All responses must strictly adhere to the JSON format.

We have defined the formalization of tool calls in § 2.4: Each tool-calling task is represented as a tuple (Q, T) , where Q is the query associated with the task, and T denotes the list of tools that the assistant can utilize. The tool-calling trajectory \mathcal{T} is a sequence of tool-response pairs $\{(a_i, r_i)\}$, which capture the interaction between the assistant's actions a and the corresponding tool responses r in the i -th step. The action a is regarded as either $(goal, tool, args)$ or $(tool, args)$ depending on whether the chain of thought (CoT) strategy is used. The test data consists of the first k steps of the tool-calling trajectory for each task, where k is randomly selected, and errors may be introduced at step k .

In an internal model-driven error task, given a tool list T , query Q , a tool-calling trajectory $\mathcal{T} = \{(a_1, r_1) \dots (a_k, r_k)\}$, and an error may be contained in a_k . The assistant is asked to generate solution $S^{pred} = (c^{pred}, \hat{a}^{pred})$ if it identifies an error in a_k , and $S^{pred} = (\hat{a}^{pred})$ otherwise. The golden solution is $S^{gt} = \{\hat{a}_1^{gt}, \hat{a}_2^{gt}\}$, where $\hat{a}_1^{gt} = a_k$ and \hat{a}_2^{gt} is the ground truth action for next subtask.

In the case of external environment error, given a tool list T , query Q , and a tool-calling trajectory $\mathcal{T} = \{(a_1, r_1) \dots (a_k, r_k)\}$, where an external error occurs in r_k . The assistant is tasked with

retrying the action a_k no more than three times, then break free from the loop and either proceed with executing the next subtasks or finish the tool call. If the predicted action $\hat{a} = a_k$, we return the erroneous response r_k to allow the assistant to proceed. Once $\hat{a} \neq a_k$ is detected, or if more than three steps are executed, we stop the assistant's reasoning and obtain a sequence of predicted solution $S^{pred} = \{\hat{a}_1^{pred}, \hat{a}_2^{pred}, \dots\}$. The golden solution is $S^{gt} = \{\hat{a}_1^{gt}, \hat{a}_2^{gt}\}$, where $\hat{a}_1^{gt} = a_k$ and \hat{a}_2^{gt} is the ground truth action for next subtask. The evaluation process is shown in the Fig. 10.

C.3.1 REFLECT

The reflect evaluator measures the model's ability to recognize the errors in tool call trajectories. For error-free trajectory where solution path is $S^{gt} = (a^{gt})$, the evaluation focuses solely on detection accuracy. If LLM predicts $S^{pred} = (a^{pred})$, the detect score is 1; otherwise, it is 0. For error-injected trajectory where solution path is $S^{gt} = (c^{gt}, a^{gt})$, the detection score is 1 if c^{pred} in prediction S^{pred} , and 0 otherwise. The evaluator then determines whether the predicted error category c^{pred} matches the ground truth c^{gt} , achieving category score 1 if the same and 0 otherwise.

C.3.2 CORRECT

The correct evaluator assesses the model's ability to correct its actions after making a mistake. For trajectories containing errors, the evaluator first verifies whether the predicted $tool^{pred}$ matches the golden answer $tool^{gt}$. If the tool prediction is correct, the tool score is 1, and the evaluator proceeds to evaluate the correctness of the input parame-

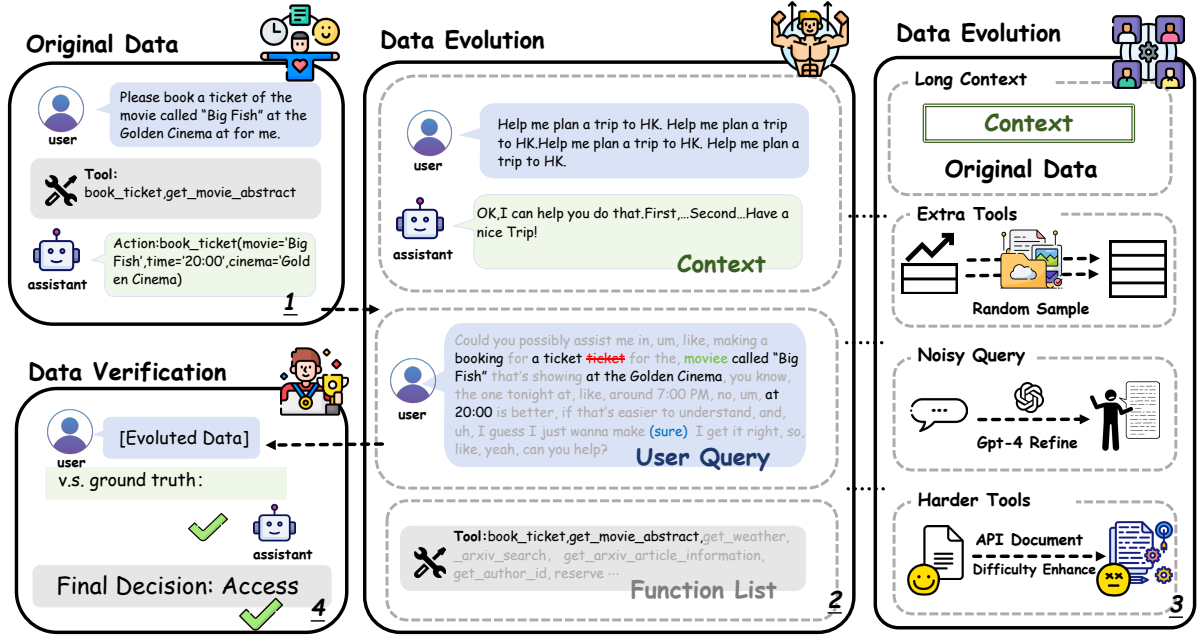


Figure 9: The framework of **Scalable and Robust Mixed Self-Evolution (SRM)**.

ters. Otherwise, both the tool and args scores are set to 0. Then, the evaluator checks whether the passed parameter keys are missing or redundant, and the args score is set to 0 if any discrepancy exists. For parameters with types such as ‘string’ or ‘any’, the evaluator uses Sentence-BERT (Reimers and Gurevych, 2019), which involves embedding the two sentences, to compute the cosine similarity between the embeddings of each predicted parameter value of $args^{pred}$ and the ground truth value $args^{gt}$ as their scores. The underlying BERT model used is all-mpnet-base-v2.² For all other parameter types, the predicted values must match the ground truth values exactly. Finally, the average score across all parameters is calculated as the args score. If the CoT strategy is applied, the evaluator uses Sentence-BERT to embed the predicted thought $thought^{pred}$ and the ground truth thought $thought^{gt}$, then calculates their cosine similarity as the thought score.

C.3.3 RETRY

The retry evaluator checks whether the predicted action \hat{a}_1^{pred} is identical to the ground truth action \hat{a}_1^{gt} , the retry score is 1 if the same and 0 otherwise.

C.3.4 SKIP

The skip evaluator first examines all predicted actions to check if there exists any $\hat{a}^{pred} \neq \hat{a}_1^{gt}$, which indicates that the model has skipped the current

retry step. If such a case \hat{a}_n^{pred} is found, the break score then set to 1. The evaluator then compares the predicted action for next subtask \hat{a}_n^{pred} with the golden answer \hat{a}_1^{gt} . The tool, args and thought score are determined using the same comparison method as in the correct evaluation.

C.3.5 FINISH

The finish evaluator first evaluates the break score in the same manner as the skip evaluator. It then checks whether the break-free action \hat{a}_n^{pred} is ‘FinishAction’. If so, the tool score is set to 1.

C.4 Experimental Details

To evaluate the pure ability of the single model, we do not use any optimization methods in the main text, such as ReAct. To assess whether the model with optimization methods exhibits a distribution comparable to the original benchmark—including indicator scores and the model’s relative strengths and weaknesses—we also generated CRITICTOOL with chain of thought (CoT). CRITICTOOL-CoT contains 810 internal model-driven error test cases and 126 external environment error test cases. Similarly, we use CRITICTOOL-CoT as the basic evolutionary dataset and obtain a total of 1,250 evolved test cases. Experimental results with CoT will be presented in the Appendix D.3.

²https://www.sbert.net/docs/pretrained_models.html

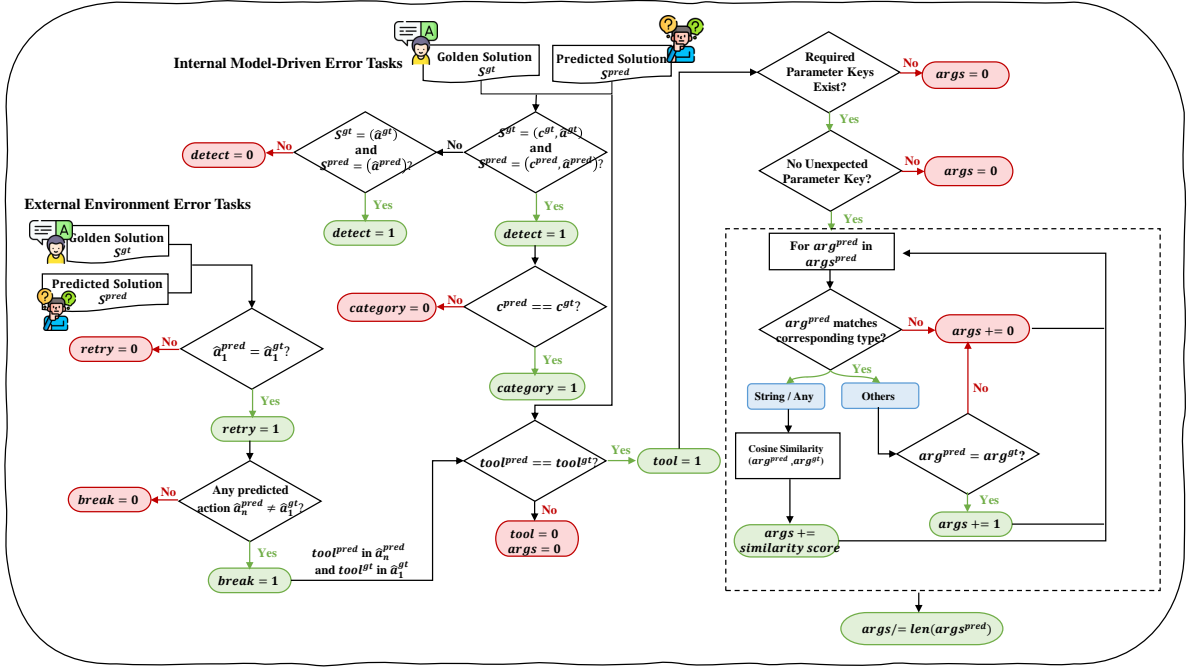


Figure 10: The framework of **Evaluation Process**.

D Additional Results

D.1 Full Results on CRITICTOOL

We show the full results on CRITICTOOL in Tab. 5, 6, 7, 8, 9 and 10.

D.2 Results of Self-Critique Performance Across Internal Error Patterns

We summarize the performance of various models across internal error patterns in our experiments, as shown in Tab. 11.

Our experimental results reveal a surprising phenomenon: even when LLMs fail to accurately identify or classify their own errors during tool calls, they are still capable of correcting these errors, which is particularly evident in tool selection errors. Although this behavior diverges from the human cognitive process, where recognizing errors typically precedes correcting them, we can still identify plausible explanations for this. During the reflection phase, LLMs heavily rely on external and explicit error signals while often overlooking the implicit errors, such as failing to obtain necessary information. This limitation stems from inadequacies in the models' instruction-following capabilities, particularly their ability to recognize subtle or implicit errors. In contrast, current training paradigms for tool use focus on enabling models to interpret the discrepancy between the expected and actual results serves as implicit feedback, allowing

models to adapt their behavior to complete tasks, even without explicitly identifying or categorizing the errors.

D.3 Full Results on CRITICTOOL-CoT

We show the full results on CRITICTOOL-CoT in Tab. 12, 13, 14, 15, 16 and 17.

E Limitation

While CRITICTOOL offers the first fine-grained and comprehensive evaluation of tool invocation self-criticism, as far as we know, it still has the following two limitations. (1) Our dataset builds upon and extends BFCL and T-eval. Despite refinement and filtering, the quality of the underlying dataset still impacts the overall quality and discriminative power of CRITICTOOL to some extent. (2) The construction of our benchmark heavily relies on GPT-4o for error generation, evolution, and verification. Although this approach ensures a high level of data quality to some extent, it may introduce potential biases inherent to the GPT-4o, potentially affecting the objectivity and robustness of the evaluation. Moreover, the dependence on high-performance LLM results in significant economic costs, posing challenges to the sustainability of large-scale benchmark development.

Future work should tackle these challenges by developing more rational and cost-effective data construction methods.

Table 5: Results of CRITICTOOL on **Base** and **Evolutionary Datasets**. **Bold** indicates the best performance across all models, while underline denotes the best performance within the same group and scale of models.

Models	Internal Model-Driven Errors								External Environment Errors								Overall	
	Reflect				Correct				Retry		Break		Skip/Finish		Args			
	Detect	Category	Tool	Args														
	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol
Close-Source Large Language Models																		
Claude3.5	85.0	78.5	60.7	50.4	87.1	81.1	80.2	74.5	45.7	34.0	57.2	55.8	22.7	22.0	26.7	26.2	57.9	53.5
GPT-3.5	73.3	78.3	61.3	66.0	72.0	67.8	58.6	50.7	12.6	9.6	92.5	87.0	54.6	49.6	46.4	41.3	62.7	59.4
GPT-4o	80.6	78.1	73.0	68.8	87.6	82.7	82.3	77.3	19.8	17.6	91.4	100.0	53.7	52.2	45.1	42.6	70.4	68.8
Open-Source Large Language Models																		
LLaMA3-8B	51.0	63.1	26.5	33.2	75.6	70.4	67.6	61.7	35.6	26.9	73.3	72.4	28.4	24.9	31.3	24.8	51.0	49.1
LLaMA3.1-8B	84.5	85.0	68.6	67.9	80.4	76.7	72.3	66.8	52.9	49.6	71.0	83.2	24.4	28.4	21.2	25.3	58.3	59.8
Qwen2.5-7B	85.1	81.7	43.1	44.4	79.6	74.2	72.1	65.3	34.2	25.7	87.6	88.6	46.0	36.9	19.7	24.7	60.3	57.4
GLM4-9B-chat	60.8	53.2	26.7	23.6	63.2	57.0	53.1	46.4	22.4	17.0	84.8	94.4	39.1	34.8	20.5	24.7	49.0	47.1
Ministral-8B	47.0	50.7	23.8	29.3	70.6	66.7	61.4	56.4	56.0	49.0	58.0	63.6	20.4	17.8	28.1	16.8	45.7	43.6
LLaMA3-70B	61.4	50.4	33.7	24.9	72.6	65.2	66.5	58.1	37.0	30.2	58.8	83.8	30.9	26.0	30.2	26.1	50.2	47.9
LLaMA3.1-70B	83.6	79.1	64.3	58.1	84.4	80.7	69.3	64.0	71.8	60.9	85.6	97.5	53.7	54.5	31.0	24.7	67.0	65.0
Qwen2.5-72B	89.4	83.4	58.9	50.7	84.5	81.9	77.9	75.5	38.8	42.0	95.1	98.4	56.9	54.6	32.4	33.3	68.8	67.1
Tool-Use-Finetuned Large Language Models																		
ToolLLaMA2-7B	0.8	0.4	0.0	0.0	4.1	2.5	0.6	0.7	1.0	0.8	1.4	0.0	0.7	1.1	0.0	0.0	1.1	0.7
ToolACE-8B	12.8	13.8	0.9	1.0	14.5	14.9	13.2	13.2	1.4	1.1	13.2	6.8	6.9	8.8	10.9	13.5	10.3	10.1
AgentLM-7B	24.9	19.3	0.0	0.0	56.0	34.1	44.1	25.5	12.1	11.3	85.1	88.5	20.4	17.9	21.0	16.0	37.1	29.8

Table 6: Results of CRITICTOOL with Only **Mixed Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	71.00	43.15	69.86	63.55	23.00	60.00	18.00	15.88	46.66
GPT-3.5	74.00	59.59	65.75	50.20	9.00	72.00	35.00	23.75	50.81
GPT-4o	81.00	70.55	74.66	67.44	15.00	100.00	44.00	33.70	63.87
Open-Source Large Language Models									
LLaMA3-8B	74.50	45.21	63.70	52.60	20.00	76.00	30.00	27.35	50.42
LLaMA3.1-8B	81.00	63.70	67.81	56.69	48.00	75.00	28.00	23.51	54.52
Qwen2.5-7B	74.50	45.21	63.70	52.60	22.00	87.00	42.00	27.35	53.97
GLM4-9B-chat	37.00	17.12	41.78	32.97	10.00	82.00	25.00	26.58	37.16
Ministral-8B	60.50	43.15	59.59	50.19	61.00	46.00	12.00	14.00	40.68
LLaMA3-70B	31.50	13.01	50.68	43.48	28.27	72.90	17.16	14.60	35.69
LLaMA3.1-70B	70.50	45.89	70.55	53.61	55.00	96.00	43.00	7.10	54.93
Qwen2.5-72B	73.50	39.73	73.97	67.63	52.00	97.00	50.00	29.92	61.70
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.50	0.00	2.05	0.77	2.27	0.00	0.00	0.00	0.59
ToolACE-8B	12.50	0.00	7.53	6.21	1.00	10.11	12.00	19.52	9.59
AgentLM-7B	7.00	0.00	13.70	9.15	9.09	81.82	2.27	3.30	17.69

Table 7: Results of CRITICTOOL with Only **Harder Tools Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	85.00	60.27	84.25	78.38	42.00	53.00	18.00	23.56	55.20
GPT-3.5	78.50	64.38	69.18	50.89	8.00	92.00	53.00	49.21	61.83
GPT-4o	88.00	82.19	86.30	82.15	22.00	100.00	55.00	41.85	69.91
Open-Source Large Language Models									
LLaMA3-8B	83.00	45.89	77.40	70.72	30.00	77.00	25.00	23.89	55.49
LLaMA3.1-8B	87.00	71.92	80.82	68.62	50.00	79.00	29.00	26.11	57.92
Qwen2.5-7B	83.00	45.89	77.40	70.72	31.32	77.01	29.60	8.25	53.90
GLM4-9B-chat	71.00	34.25	64.38	52.06	22.00	100.00	44.00	29.04	55.05
Ministral-8B	52.50	32.88	68.49	58.84	18.00	92.00	12.00	5.15	44.91
LLaMA3-70B	67.50	35.62	73.29	65.19	36.00	87.00	31.00	22.59	53.97
LLaMA3.1-70B	88.00	67.12	83.56	70.57	71.55	94.54	44.25	4.30	63.67
Qwen2.5-72B	87.00	52.05	84.25	79.52	53.00	100.00	60.00	40.77	71.24
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.50	0.00	0.00	0.00	1.61	0.00	0.00	0.00	0.13
ToolACE-8B	17.50	0.00	21.23	17.84	0.29	0.00	0.00	0.00	6.63
AgentLM-7B	23.50	0.00	43.84	30.18	10.26	92.31	26.92	27.76	36.01

Table 8: Results of CRITICTOOL with With Only **Noisy Query Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	79.00	48.63	80.82	75.08	36.00	48.00	17.00	20.15	50.72
GPT-3.5	77.50	65.75	64.38	46.30	15.00	92.00	52.00	43.57	59.81
GPT-4o	77.50	69.18	80.82	76.83	20.00	100.00	54.00	44.24	69.05
Open-Source Large Language Models									
LLaMA3-8B	45.50	20.55	71.23	64.08	40.00	78.00	28.00	28.05	49.01
LLaMA3.1-8B	87.00	71.23	77.40	70.33	54.00	87.00	35.00	31.93	57.77
Qwen2.5-7B	84.50	39.04	71.92	62.99	27.00	99.00	38.00	26.17	58.42
GLM4-9B-chat	57.00	22.60	54.79	45.44	26.00	100.00	37.00	22.20	48.17
Ministral-8B	46.00	23.97	65.75	55.94	56.00	58.00	21.00	23.37	43.41
LLaMA3-70B	55.00	24.66	65.07	60.25	43.18	82.55	32.89	34.84	51.46
LLaMA3.1-70B	84.00	67.12	78.77	64.45	79.00	100.00	71.00	38.23	71.93
Qwen2.5-72B	88.50	58.90	79.45	73.76	52.00	99.00	58.00	30.16	68.40
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.00	0.00	2.74	0.34	0.00	0.00	2.35	0.00	0.82
ToolACE-8B	12.00	2.74	13.01	12.84	3.00	0.00	7.00	7.18	7.33
AgentLM-7B	26.50	0.00	46.58	35.57	8.25	91.75	22.45	22.67	35.91

Table 9: Results of CRITICTOOL with Only **Extra Tools Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	81.50	56.16	82.88	75.02	42.00	54.00	24.00	33.30	56.25
GPT-3.5	80.00	69.18	71.23	52.45	8.00	83.00	53.00	53.48	62.29
GPT-4o	70.50	59.59	85.62	78.51	17.00	100.00	55.00	44.34	68.38
Open-Source Large Language Models									
LLaMA3-8B	82.50	46.58	77.40	66.95	23.47	67.35	22.45	31.15	53.87
LLaMA3.1-8B	86.50	70.55	78.77	68.26	43.00	88.00	22.00	21.17	57.59
Qwen2.5-7B	82.00	45.89	77.40	68.58	23.56	81.32	33.91	30.35	57.70
GLM4-9B-chat	53.00	25.34	63.70	52.49	12.00	90.00	39.00	27.98	49.41
Ministral-8B	49.50	26.03	69.86	57.87	57.00	57.00	17.00	14.78	42.88
LLaMA3-70B	62.00	35.62	68.49	60.04	31.91	80.85	29.79	35.25	49.52
LLaMA3.1-70B	79.50	59.59	82.88	62.01	62.00	97.00	59.00	40.48	68.21
Qwen2.5-72B	87.50	54.11	86.30	76.90	32.00	97.00	57.00	34.81	68.56
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.50	0.00	6.85	2.40	0.00	0.00	3.03	0.00	1.89
ToolACE-8B	19.50	1.37	20.55	17.52	0.00	6.00	15.00	24.59	13.64
AgentLM-7B	26.00	0.00	42.47	34.42	17.14	88.57	14.29	13.80	32.49

Table 10: Results of CRITICTOOL with Only **Long Context Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	76.00	43.84	87.67	80.23	27.00	64.00	33.00	38.00	58.77
GPT-3.5	81.50	71.23	68.49	53.67	8.00	96.00	55.00	36.57	62.13
GPT-4o	73.50	62.33	86.30	81.54	14.00	100.00	53.00	48.76	69.72
Open-Source Large Language Models									
LLaMA3-8B	30.00	7.53	62.33	54.34	21.28	63.83	19.15	13.51	36.79
LLaMA3.1-8B	83.50	62.33	78.77	70.29	53.00	87.00	28.00	23.95	55.43
Qwen2.5-7B	84.50	45.89	80.82	71.63	24.71	98.85	41.09	31.25	56.82
GLM4-9B-chat	48.00	18.49	60.27	49.12	15.00	100.00	29.00	17.91	45.85
Ministral-8B	45.00	20.55	69.86	59.39	53.00	65.00	27.00	26.51	46.37
LLaMA3-70B	36.00	15.75	68.49	61.74	11.89	95.65	19.39	23.42	46.07
LLaMA3.1-70B	73.50	50.68	87.67	69.18	37.00	100.00	55.00	33.55	66.08
Qwen2.5-72B	80.50	48.63	85.62	79.75	21.00	99.00	48.00	30.65	65.42
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.50	0.00	0.68	0.00	0.00	0.00	0.00	0.00	0.15
ToolACE-8B	7.50	0.68	12.33	11.64	1.00	0.00	10.00	16.19	8.39
AgentLM-7B	13.50	0.00	23.97	18.33	11.81	88.19	23.61	12.50	26.93

Table 11: **Self-Critique Evaluation** on different error patterns.

Models	Tool Sel. Errors		Tool Halluc. Errors		Param. Key Errors		Param. Value Errors	
	Reflect	Correct	Reflect	Correct	Reflect	Correct	Reflect	Correct
Close-Source Large Language Models								
Claude3.5	10.15	56.29	93.29	65.74	93.21	90.59	94.11	90.80
GPT-3.5	7.32	32.81	80.10	27.89	82.65	79.07	86.96	66.28
GPT-4o	23.42	59.18	97.72	70.43	79.65	92.81	86.17	90.22
Open-Source Large Language Models								
LLaMA3-8B	7.68	41.58	70.30	52.29	61.39	83.07	67.79	78.12
LLaMA3.1-8B	19.48	41.29	97.49	54.69	98.47	88.90	92.60	82.60
Qwen2.5-7B	28.14	37.61	96.51	57.68	97.40	85.96	93.38	85.25
GLM4-9B-chat	9.58	18.35	61.42	42.34	55.98	69.83	62.93	55.86
Ministral-8B	4.27	34.42	70.07	42.38	23.68	77.86	29.43	70.35
LLaMA3-70B	8.15	43.09	70.21	55.33	57.48	76.95	54.99	66.00
LLaMA3.1-70B	14.11	49.66	94.51	51.17	90.79	78.61	91.53	83.18
Qwen2.5-72B	36.92	55.91	94.03	59.34	95.37	91.08	97.03	93.73
Tool-Use-Finetuned Large Language Models								
ToolLLaMA2-7B	0.29	0.00	0.76	0.00	0.30	0.93	1.00	1.65
ToolACE-8B	0.28	11.11	3.25	5.01	2.74	19.16	4.31	13.48
AgentLM-7B	0.56	20.70	1.26	22.83	0.30	50.62	0.68	40.53

Table 12: Results of CRITICTOOL-CoT on **Base** and **Evolutionary Datasets**.

Models	Internal Model-Driven Errors								External Environment Errors								Overall	
	Reflect				Correct				Retry				Skip/Finish				Args	
	Detect		Category		Tool		Args		Break		Tool		Args		Base		Evol	
	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol	Base	Evol
Close-Source Large Language Models																		
Claude3.5	91.7	82.5	71.2	55.9	90.7	85.5	83.8	77.5	37.3	25.2	94.4	66.4	36.9	23.6	51.4	35.4	71.8	58.4
GPT-3.5	67.0	78.8	52.1	66.3	84.4	75.1	70.3	57.9	15.1	7.8	81.0	78.8	<u>63.5</u>	<u>55.4</u>	48.5	42.5	64.8	61.4
GPT-4o	91.4	86.9	86.5	81.2	90.4	85.3	85.1	79.3	<u>45.6</u>	<u>39.1</u>	100.0	100.0	47.6	45.7	62.9	60.7	78.0	74.4
Open-Source Large Language Models																		
LLaMA3-8B	70.9	69.6	48.9	38.9	79.8	76.3	74.0	68.7	43.7	40.9	82.9	74.6	55.6	38.7	29.9	26.6	62.5	55.6
LLaMA3.1-8B	90.2	<u>84.5</u>	<u>77.7</u>	<u>72.7</u>	<u>85.3</u>	80.5	<u>79.1</u>	71.4	52.0	49.4	89.3	93.6	<u>56.3</u>	<u>53.4</u>	28.3	29.3	<u>70.1</u>	<u>67.4</u>
Qwen2.5-7B	88.5	80.7	49.1	43.3	83.5	<u>81.1</u>	77.2	<u>73.8</u>	79.4	68.9	<u>92.1</u>	<u>96.5</u>	56.0	50.3	34.9	27.3	69.3	65.2
GLM4-9B-chat	78.4	58.4	33.0	27.0	76.5	68.5	65.2	56.7	28.2	20.5	86.1	92.0	49.6	42.1	<u>42.0</u>	<u>36.3</u>	60.4	53.9
Ministral-8B	45.6	44.1	20.5	21.6	76.1	73.3	68.7	64.0	69.0	58.4	40.5	53.2	15.5	13.0	23.6	14.3	43.7	42.2
LLaMA3-70B	69.1	56.4	42.8	34.4	83.3	71.1	75.8	62.7	56.4	40.1	83.2	87.9	50.0	43.5	25.4	26.9	61.7	54.9
LLaMA3.1-70B	90.0	78.9	<u>75.8</u>	<u>60.8</u>	<u>85.8</u>	81.0	73.4	67.7	70.2	64.5	96.4	98.4	65.9	57.2	36.8	26.1	73.8	66.7
Qwen2.5-72B	<u>91.7</u>	<u>81.8</u>	57.9	46.8	85.3	<u>81.8</u>	<u>79.6</u>	<u>75.2</u>	69.8	<u>65.6</u>	<u>96.8</u>	<u>98.4</u>	68.3	61.2	57.4	49.0	<u>76.6</u>	<u>71.0</u>
Tool-Use-Finetuned Large Language Models																		
ToolLLaMA2-7B	0.4	0.6	0.0	0.0	0.9	1.5	0.2	0.2	0.0	1.5	0.4	1.2	0.0	0.0	0.0	0.0	0.3	0.6
ToolACE-8B	14.6	11.1	1.8	1.0	20.4	19.5	18.2	16.9	4.0	2.2	10.7	2.4	7.1	10.2	10.5	14.8	11.9	10.9
AgentLM-7B	<u>25.2</u>	<u>15.6</u>	0.0	0.0	<u>48.6</u>	<u>31.1</u>	<u>35.4</u>	<u>21.7</u>	<u>47.5</u>	<u>41.6</u>	<u>48.3</u>	2.0	<u>19.4</u>	<u>14.8</u>	<u>16.4</u>	<u>15.0</u>	<u>30.1</u>	<u>16.3</u>

Table 13: Results of CRITICTOOL-CoT with **Only Mixed Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	78.00	52.74	76.03	67.67	15.00	62.00	24.00	27.52	52.41
GPT-3.5	70.50	54.79	69.18	58.27	3.00	78.00	39.00	27.63	53.49
GPT-4o	84.50	78.08	74.66	70.17	33.31	100.00	42.49	41.66	67.27
Open-Source Large Language Models									
LLaMA3-8B	73.50	48.63	74.66	69.94	42.58	84.62	26.54	24.16	56.33
LLaMA3.1-8B	81.50	70.55	73.29	64.19	51.00	88.00	43.00	22.77	61.44
Qwen2.5-7B	73.50	48.63	74.66	69.94	57.00	93.00	39.00	24.16	60.18
GLM4-9B-chat	38.00	15.07	54.11	44.43	12.00	84.00	30.00	24.24	41.42
Ministral-8B	52.50	33.56	67.12	57.26	75.00	30.00	4.00	2.00	36.41
LLaMA3-70B	46.00	30.82	59.59	51.27	23.96	78.12	22.92	18.67	43.47
LLaMA3.1-70B	71.50	55.48	71.92	59.48	63.00	97.00	56.00	17.24	61.09
Qwen2.5-72B	77.50	47.26	76.03	69.25	66.00	97.00	49.00	30.96	64.11
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.00	0.00	0.00	0.00	4.20	3.93	0.00	0.00	0.80
ToolACE-8B	9.00	0.00	20.55	16.38	2.00	0.81	8.00	10.41	9.42
AgentLM-7B	8.00	0.00	10.27	6.44	27.45	0.84	13.64	18.34	9.60

Table 14: Results of CRITICTOOL-CoT with **Only Harder Tools Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	85.50	60.96	86.99	78.54	30.00	69.00	27.00	43.80	61.94
GPT-3.5	81.50	69.18	77.40	64.30	11.00	82.00	63.00	45.71	65.48
GPT-4o	90.00	85.62	89.04	82.84	39.00	100.00	45.00	68.61	77.34
Open-Source Large Language Models									
LLaMA3-8B	84.50	41.78	83.56	75.80	43.00	84.00	58.00	41.22	66.17
LLaMA3.1-8B	86.50	73.97	82.88	71.60	45.00	96.00	57.00	31.97	69.21
Qwen2.5-7B	84.50	41.78	83.56	75.80	74.60	94.44	46.83	5.92	62.34
GLM4-9B-chat	69.50	34.93	79.45	65.87	31.39	83.84	46.51	34.75	58.58
Ministral-8B	46.00	23.97	73.29	62.28	36.00	96.00	9.00	13.92	46.97
LLaMA3-70B	72.50	45.21	80.82	69.43	54.00	90.00	54.00	29.69	63.06
LLaMA3.1-70B	87.00	67.81	86.30	70.90	70.63	98.02	48.02	6.38	65.45
Qwen2.5-72B	87.50	47.95	84.93	77.81	67.00	98.00	61.00	49.16	72.53
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.50	0.00	2.05	0.50	1.02	3.38	0.00	0.00	0.99
ToolACE-8B	14.50	2.05	21.23	17.59	0.00	0.08	10.00	18.65	11.79
AgentLM-7B	21.50	0.00	50.68	36.03	50.00	3.07	13.27	15.06	22.37

Table 15: Results of CRITICTOOL-CoT with With Only Noisy Query Evolution Data.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	86.50	60.27	88.36	81.74	21.00	70.00	21.00	33.91	59.98
GPT-3.5	80.50	69.18	77.40	55.04	13.00	81.00	54.00	42.99	62.18
GPT-4o	88.00	82.88	88.36	83.05	45.00	100.00	46.00	64.18	76.58
Open-Source Large Language Models									
LLaMA3-8B	68.00	44.52	76.03	69.41	52.00	82.00	49.00	24.94	59.06
LLaMA3.1-8B	86.50	75.34	84.93	76.83	58.00	97.00	57.00	25.60	70.29
Qwen2.5-7B	84.50	40.41	84.93	79.15	83.00	98.00	55.00	35.56	69.54
GLM4-9B-chat	71.50	32.19	78.08	66.10	24.00	100.00	45.00	40.57	61.03
Ministral-8B	37.50	14.38	77.40	69.18	59.00	50.00	13.00	20.00	42.58
LLaMA3-70B	65.50	38.36	75.34	67.55	61.00	89.00	57.00	33.24	61.76
LLaMA3.1-70B	86.00	73.29	83.56	71.29	79.00	100.00	61.00	30.15	71.78
Qwen2.5-72B	85.00	50.00	86.99	81.75	77.00	99.00	66.00	52.21	75.24
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	1.50	0.00	2.74	0.68	1.00	2.69	0.00	0.00	1.12
ToolACE-8B	12.50	2.74	21.92	18.92	5.00	4.92	12.00	16.90	12.97
AgentLM-7B	23.00	0.00	45.21	33.44	43.00	3.79	14.00	14.30	21.06

Table 16: Results of CRITICTOOL-CoT with Only Extra Tools Evolution Data.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	85.50	62.33	89.73	82.50	30.00	63.00	20.00	32.64	59.46
GPT-3.5	80.00	67.81	76.71	57.88	6.00	71.00	62.00	47.31	62.32
GPT-4o	86.50	80.82	87.67	80.69	41.00	100.00	46.00	62.99	75.38
Open-Source Large Language Models									
LLaMA3-8B	82.50	40.41	80.14	72.61	35.71	68.37	31.63	28.72	56.30
LLaMA3.1-8B	84.50	73.97	81.51	70.50	44.00	94.00	55.00	30.33	67.75
Qwen2.5-7B	83.50	41.10	82.19	71.53	67.06	98.02	56.75	37.07	67.65
GLM4-9B-chat	67.00	32.19	66.44	52.22	20.00	92.00	51.00	49.14	57.54
Ministral-8B	44.00	21.23	75.34	65.33	68.00	38.00	16.00	17.75	41.79
LLaMA3-70B	61.50	37.67	78.77	68.47	43.88	88.78	52.04	33.85	60.40
LLaMA3.1-70B	83.00	64.38	83.56	66.37	62.00	97.00	63.00	34.60	69.52
Qwen2.5-72B	85.50	52.05	83.56	75.42	62.00	98.00	68.00	61.86	74.88
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	1.00	0.00	2.05	0.00	1.09	4.08	0.00	0.00	1.07
ToolACE-8B	10.00	0.00	20.55	19.57	3.00	3.17	12.00	15.61	11.79
AgentLM-7B	22.00	0.00	43.84	31.44	40.22	1.23	20.65	27.16	22.86

Table 17: Results of CRITICTOOL-CoT with Only **Long Context Evolution Data**.

Models	Internal Model-Driven Errors				External Environment Errors				Overall
	Reflect		Correct		Retry	Break	Skip/Finish		
	Detect	Category	Tool	Args			Tool	Args	
Close-Source Large Language Models									
Claude3.5	77.00	43.15	86.30	77.29	30.00	68.00	26.00	39.37	58.06
GPT-3.5	81.50	70.55	74.66	54.04	6.00	82.00	59.00	48.87	63.29
GPT-4o	85.50	78.77	86.99	79.96	37.00	100.00	49.00	66.24	75.61
Open-Source Large Language Models									
LLaMA3-8B	39.50	19.18	67.12	55.96	31.25	54.17	28.12	14.03	40.34
LLaMA3.1-8B	83.50	69.86	80.14	73.67	49.00	93.00	55.00	35.84	68.43
Qwen2.5-7B	77.50	44.52	80.14	72.53	62.70	99.21	53.97	33.87	66.29
GLM4-9B-chat	46.00	20.55	64.38	55.09	15.00	100.00	38.00	32.87	50.96
Ministral-8B	40.50	15.07	73.29	65.85	54.00	52.00	23.00	17.95	43.07
LLaMA3-70B	36.50	19.86	60.96	56.59	17.86	93.68	31.63	19.00	45.81
LLaMA3.1-70B	67.00	43.15	79.45	70.28	48.00	100.00	58.00	42.10	65.89
Qwen2.5-72B	73.50	36.99	77.40	71.57	56.00	100.00	62.00	51.04	68.15
Tool-Use-Finetuned Large Language Models									
ToolLLaMA2-7B	0.00	0.00	0.68	0.00	0.00	3.32	0.00	0.00	0.60
ToolACE-8B	9.50	0.00	13.01	12.13	1.00	0.24	9.00	12.25	8.00
AgentLM-7B	3.50	0.00	5.48	1.03	47.32	4.00	12.50	0.00	6.17

Standardization of Test Data

System Prompt:

You are a assistant who can utilize external tools. You can call the following tools:

[API List]

To use a tool, please directly generate the response in JSON format. Do NOT add comments (//) in the response.

```
{
  "name": string,    // tool name to call
  "args": Record<string, any>  // input params required by current tool in JSON format
}
```

If you already know the answer, or you do not need to use tools, please using the following format to reply:

```
{
  "name": "FinishAction",
  "args": "string"    // your final response to the query
}
```

Remember: only generate ONE step each time.

User: Help me maintain a social media presence by crafting a tweet that states, 'Managed to archive important data files!' using the hashtags #DataManagement and #Efficiency.

Assistant:

```
{
  "name": "TwitterAPI.authenticate_twitter",
  "args": {"username": "dr_smith", "password": "securePass123"}
}
```

Function:

```
{
  "authentication_status": true
}
```

Prediction: The assistant's tool-calling action for the next step.

Ground True:

```
{
  "name": "TwitterAPI.post_tweet",
  "args": {"content": "Managed to archive important data files!", "tags": ["#DataManagement", "#Efficiency"]}
}
```

Figure 11: An example of standardized test data.

Insight into Internal Model-Driven Error

User: Retrieve information about author John Smith.

Ground True:

```
{
  "name": "GoogleScholar.get_author_id",
  "args": {"mauthors": "John Smith"}
}
```

Error Prediction:

Assistant

```
{
  "name": "GoogleScholar.get_author_information",
  "args": {"author_id": "John Smith"}
}
```

Calls API

Function

```
{
  "text": "Fully empty",
  "error": "GoogleScholar hasn't returned any results for this author_id."
}
```

Continue Prediction:

Assistant

```
{
  "name": "GoogleScholar.get_author_id",
  "args": {"mauthors": "John Smith"}
}
```

Calls API

Function

```
{ "author_id": "0juNTSgAAAAJ" }
```

Correct

Assistant

```
{
  "name": "GoogleScholar.get_author_information",
  "args": {"author_id": "0juNTSgAAAAJ"}
}
```

Assistant

```
{
  "name": "FinishAction",
  "args": "I am unable to find any information about author John Smith in GoogleScholar."
}
```

Figure 12: Examples of LLMs' performance in the presence of internal errors.

Insight into External Environment Error

User: Read the content of "/home/user/Documents/file.txt", then find the upcoming action movies in China.

Prediction:

Assistant

```
{
  "name": "FileOperation.read_file",
  "args": {"file_path": "/home/user/Documents/file.txt"}
}
```

Calls API

Function

```
{
  "text": "Error: [Errno 13] Permission denied: '/home/user/Documents/file.txt'"
}
```

Continue Prediction:

Assistant

```
{
  "name": "FileOperation.read_file",
  "args": {"file_path": "/home/user/Documents/file.txt"}
}
```

Calls API

Function

```
{
  "text": "Error: [Errno 13] Permission denied: '/home/user/Documents/file.txt'"
}
```

↓

Assistant

```
{
  "name": "FilmDouban.coming_out_filter",
  "args": {}
}
```

Assistant

```
{
  "name": "FileOperation.read_file",
  "args": {"file_path": "/home/user/Documents/file.txt"}
}
```

Calls API

Function

```
{
  "text": "Error: [Errno 13] Permission denied: '/home/user/Documents/file.txt'"
}
```

⋮ (Endlessly Retry)

Assistant

```
{
  "name": "FileOperation.read_file",
  "args": {"file_path": "/home/user/Documents/file.txt"}
}
```

Figure 13: Examples of LLMs' performance in the presence of external errors.

Collected Tool-Calling Trajectory

System Prompt:

Same as standardized test data.

User: Tool-use task.

Assistant:

```
{
  "name": "Golden answer API1",
  "args": "Golden answer parameters1"
}
```

Function:

```
{Response from API1 to the input parameters1}
```

Assistant:

```
{
  "name": "Golden answer API2",
  "args": "Golden answer parameters2"
}
```

Function:

```
{Response from API2 to the input parameters2}
```

⋮

Assistant:

```
{
  "name": "FinishAction",
  "args": "The answer of the task is ..."
}
```

Figure 14: An example of collected tool-calling trajectories.

Refined API Documentation

```
{
  "name": "TravelAPI.cancel_booking",
  "description": "Cancel a booking",
  "required_parameters": [
    {
      "name": "access_token",
      "type": "string",
      "description": "[Required] The access token obtained from the authenticate"
    },
    {
      "name": "booking_id",
      "type": "string",
      "description": "[Required] The ID of the booking"
    }
  ],
  "optional_parameters": [],
  "return_data": [
    {
      "name": "cancel_status",
      "description": "The status of the cancellation, True if successful, False if failed"
    },
    {
      "name": "error",
      "description": "The error message if the cancellation failed"
    }
  ]
}
```

Figure 15: An example refined API documentation: TravelAPI.

System Prompt:**Character Introduction**

You are a large language modeling engineer, and your current task is to modify some conversation datas of large language model interacting with some external tool APIs. Your goal is to modify the content of the last reply of assistant in the correct dialog so that an error occurs and matches the error category I have given.

Description of the Dialogues Structure

- User presents the task and describes the problems to be solved.
- Assistant replies to solve the problems, may call the tool API or give the answer directly.
- Function is a tool API return that provides actual datas or the results of performing a specific action.
- The interaction consists of several steps, and the assistant solves the problems step-by-step by calling functions.

Your Task

- Find the dialog to be modified: identify the last assistant response in each dialog that is the target of the message you need to modify.
- Understanding error categories: I will provide you with a specific error category, and you need to analyze the original dialog according to the error category and find out what needs to be modified, making sure that each step of your analysis is clear and reasonable.
- Conduct modifications: make the appropriate modifications based on the error category so that the dialog contains errors that match that error category.

Response Format

Follow the JSON format to output only the modified dialog without redundancy, and do not add comments (//) in the response.

```
{
  "role": "assistant",
  "content": "{('thought': string, // goal at current step)
               'name': string, // tool name to call
               'args': Record<string, any>} // input params required by current tool in JSON
            format"
}
```

Notes

- Accuracy of JSON format: Please strictly follow the reply format, and output only the modified wrong tool call action of assistant.
- Reasonability of tool call: even if the error is generated, the called tool and its argument settings should be within a reasonable range, and the error should have some relevance to the correct dialog.
- Keep the chain of thought clear: although it is a simulation of the dialog and errors, assistant's thought process still needs to be clear and reasonable. Even if an error occurs, the logic of the assistant's reasoning when calling the tool should be complete.

Modification Example

[Randomly select 3 instances of a specific pattern of error from benchmark tests as few-shot.]

User:

Now I'll provide you with the error type and the correct dialog trajectory, please modify the last assistant's response to correspond to the error type.

Error Type: Tool Select Error/Tool Hallucination Error/Parameters Key Error/Parameters Value Error
Correct Dialog Trajectory: [randomly select the first k steps of tool call trajectory]

Figure 16: An example prompt of Error Diversification.

System Prompt:

Imagine you are an API Server operating within a specialized tool, which contains a collection of distinct APIs. Your role is to deeply understand the function of each API based on their descriptions in the API documentation. As you receive specific inputs for individual API calls within this tool, analyze these inputs to determine their intended purpose. Your task is to craft a response that aligns with the expected output of the API, guided by the provided examples.

Please note that your answer should not contain anything other than a json format object, which should be parsable directly to json, which is as follows:

```
{
  "error": "",
  "response": "<Your_Response>"
}
```

The error field should returns an explicit error message describing the cause of the error if there are any errors in the API Input. The response field must adhere strictly JSON format. <Your_Response> should contain the return_data you formulate based on the API's functionality and the input provided. Ensure that your responses are meaningful, directly addressing the API's intended functionality.

API calls may fail for various reasons, such as invalid input parameters, authentication issues, or server errors. Your goal is to generate a response that accurately reflects the API's intended functionality, even if the input parameters are incorrect. Your response should be informative and relevant to the API's purpose, providing a clear and concise explanation of the expected output based on the input provided. If the user explicitly requests messages about failed api calls, and most of the examples provided get an error response despite passing in correct and valid parameters, please generate a failed tool call response containing some external environment errors. The external environment errors include rate limit exceeded, permission denied, maximum quota exceeded, timeout, connection error and so on. Please randomly select one kind of error above, the error message should match the corresponding api as much as possible, and don't show the words "external environment error".

Note that:

- You should strictly validate the parameters of the API Input to ensure all required_parameters are provided, the value of each parameter strictly conforms to the type specified in the api documentation, and there are no redundant parameter keys passed in. Be careful to identify the types of incoming parameters, even if they are the same as those specified by required_parameters when converted to strings, a different type can cause an error.
- If there is no error in the API Input and no explicit require by user, you should fill in the response field according to the rules, and the error field should remain empty. Otherwise, you should fill in the error field according to the rules, and the response field should remain empty.
- The response and error fields are not allowed to be filled in at the same time, you are only allowed to fill in one depending on the situation.
- Your response should be around 100 to 200 words, containing rich information given the api input parameters. Keep Your answer short and simple.

User:

API Documentation:

```
{api_doc}
```

API Examples:

```
{api_cache}
```

API Input:

```
{input args}
```

Figure 17: Prompt of API simulator.

Noisy Query Evolution

System Prompt:

Your Task

- You are a helpful assistant and will receive a request from a user. This request is sent to a task related to the LLM model.
- Your task is to make this request as complex as possible, such as adding irrelevant information, adding confusing concepts that are irrelevant to the final task, add typos that do not affect the task, etc.

Response Format

Please follow the JSON format and output according to the following structure

```
{
  "Query": string,    // the refined query
  "Explanation": string, // the reason why you refine the query
}
```

Remember: be careful NOT to affect the completion of the task.

User: Here is the user query to be refined: Copy the txt contents of the 'Quarter1_Reports' directory and place it in a new directory naming it 'Archived_Quarter1'.

Figure 18: An example prompt of Noisy Query Evolution.

Harder Tools Evolution

System Prompt:

Your Task

- You are a helpful expert. You will receive an API document. You need to change the description of this api but do not change other parts, especially parameters, etc.
- You can change the expression to make it more verbose. Do not change the original meaning of the description.

Response Format

Please follow the JSON format and output according to the following structure

```
{
  "API Document": dict,    // the refined API document
  "Explanation": string,   // the reason why you refine the API document
}
```

Remember: be careful NOT to affect the completion of the API.

User: Here is the API document to be refined:

```
{
  "name": "TimeTool.get_curr_time",
  "description": "Retrieve the current date and time",
  "required_parameters": [],
  "optional_parameters": [],
  "return_data": [
    {
      "name": "time",
      "description": "The current date and time in the format YYYY-MM-DD HH:MM"
    }
  ]
},
```

Figure 19: An example prompt of Harder Tools Evolution.

The verification of Long Context

System Prompt:

Your Task

- You are a helpful expert. You will receive a context from LLM and a user query task. Please judge whether the context will affect the task.
- Please be strict on this question. If it will affect, please reply Yes. If it will not affect, please reply No.

Response Format

Please follow the JSON format and output according to the following structure

```
{
  "Result": string,    // Yes or No
  "Reason": string,    // the reason why you think the context will or will not affect the task
}
```

User: Here is the context:

```
{
  "role": "user",
  "content": "...",
},
{
  "role": "assistant",
  "content": "... the context extracted from LongBench
}
```

and the user task is:

I am planning a trip from Times Square to Central Park in New York City. I'd like to know the best path to take, such as walking, biking, or taking public transportation.

Figure 20: An example prompt of the verification of Long Context.

The verification of Noisy Query

System Prompt:

Your Task

- You are a helpful expert. You will receive two user queries: A and B. You need to determine whether B completely contains the tasks in A and whether there is no ambiguity and typo in the important expression parts.
- If there is no ambiguity, output Yes, and if there is ambiguity, output No.

Response Format

Please follow the JSON format and output according to the following structure

```
{
  "Result": string,    // Yes or No
  "Reason": string,    // the reason why there is or is not ambiguity
}
```

User: Here is the user query A:

I am planning a trip from Times Square to Central Park in New York City. I'd like to know the best path to take, such as walking, biking, or taking public transportation. // the origin user query

Here is the user query B:

I am in the process of meticulously planning an excursion from the bustling Times Square to the serene Central Park in the heart of New York City. I am quite curious to discover the most optimal route to embark upon for this journey, whether it be the leisurely stroll of walking, the environmentally friendly and energetic biking, or the efficient and convenient public transportation system. Each option presents its own unique set of advantages and challenges, and I am eager to weigh them all carefully. // the new evolved user query

Figure 21: An example prompt of the verification of Noisy Query.

The verification of Extral Tools

System Prompt:

Your Task

- You are a helpful expert. You will receive two tool lists: tool list A and B. Your task is to determine whether there are particularly similar functions in these two function lists.
- If they are particularly similar, reply yes, otherwise reply no. Please be strict on this question.

Response Format

Please follow the JSON format and output according to the following structure

```
{
  "Result": string,    // Yes or No
  "Reason": string,    // the reason why the two tool lists are similar or different
}
```

User: Here is the tool list A:

```
{
  "name": "Tool 1",
  "description": "...",
  "required_parameters": [],
  "optional_parameters": [],
  "return_data": [
    "..."
  ]
},
{
  "name": "Tool 2",
  "description": "...",
  "required_parameters": [],
  "optional_parameters": [],
  "return_data": [
    "..."
  ]
},
// the original tool list
```

Here is the tool list B:

```
{
  "name": "Tool 3",
  "description": "...",
  "required_parameters": [],
  "optional_parameters": [],
  "return_data": [
    "..."
  ]
},
{
  "name": "Tool 4",
  "description": "...",
  "required_parameters": [],
  "optional_parameters": [],
  "return_data": [
    "..."
  ]
},
// the new added tool list
```

Figure 22: An example prompt of the verification of Extra Tools.

The verification of Harder Tools

System Prompt:

Your Task

- You will receive two API documents: API A and B. Your task is to determine whether the two APIs are equivalent, that is, whether the corresponding functions have the same parameters and whether the descriptions have the same meaning.
- The expressions may be slightly different, ignore typos).
- If they are equivalent, answer Yes, otherwise answer No.

Response Format

Please follow the JSON format and output according to the following structure

```
{
  "Result": string,    // Yes or No
  "Reason": string,    // the reason why the two APIs are equivalent or different
}
```

User: Here is the API A:

```
{
  "name": "Tool 1",
  "description": "...",
  "required_parameters": [],
  "optional_parameters": [],
  "return_data": [
    "..."
  ]
}, //the origin API document
```

Here is the API B:

```
{
  "name": "Tool 2",
  "description": "...",
  "required_parameters": [],
  "optional_parameters": [],
  "return_data": [
    "..."
  ]
}, //the new evolved API document
```

Figure 23: An example prompt of the verification of Harder Tools.

An example of BFCL

System Prompt:

Your Task

- You are a helpful expert. You will receive a context from LLM and a user query task. Please judge whether the context will affect the task.
- Please be strict on this question. If it will affect, please reply Yes. If it will not affect, please reply No.

Response Format

Please follow the JSON format and output according to the following structure

```
{
  "Result": string,    // Yes or No
  "Reason": string,    // the reason why you think the context will or will not affect the task
}
```

User: Here is the context:

```
{
  "role": "user",
  "content": "...",
},
{
  "role": "assistant",
  "content": "... the context extracted from LongBench
}
```

and the user task is:

```
I am planning a trip from Times Square to Central Park in New York City. I'd like to know the best
path to take, such as walking, biking, or taking public transportation.
```

Figure 24: An example prompt of the verification of Long Context.