#### **ACEBench: A Comprehensive Evaluation of LLM Tool Usage**

Chen Chen<sup>1,2†</sup>, Xinlong Hao<sup>3†</sup>, Weiwen Liu<sup>4\*</sup>, Xu Huang<sup>1</sup>, Xingshan Zeng<sup>3</sup>, Shuai Yu<sup>3</sup>, Dexun Li<sup>3</sup>, Yuefeng Huang<sup>1</sup>, Xiancgcheng Liu<sup>1,2</sup>, Xinzhi Wang<sup>3</sup>, Wu Liu<sup>1\*</sup>,

<sup>1</sup>University of Science and Technology of China, <sup>2</sup>Shanghai Innovation Institute, <sup>3</sup>Huawei Noah's Ark Lab, <sup>4</sup>Shanghai Jiao Tong University,

chenchen0318@mail.ustc.edu.cn haoxinlong@huawei.com

#### **Abstract**

Large Language Models (LLMs) have demonstrated significant potential in decision-making and reasoning, particularly when integrated with various tools to effectively solve complex problems. However, existing benchmarks for evaluating LLMs' tool usage face several limitations: (1) limited evaluation scenarios, often lacking assessments in real multi-turn dialogue contexts; (2) narrow evaluation dimensions, with insufficient detailed assessments of how LLMs use tools; and (3) reliance on LLMs or real API executions for evaluation, which introduces significant overhead. To address these challenges, we introduce ACEBench, a comprehensive benchmark for assessing tool usage in LLMs. ACEBench categorizes data into three primary types based on evaluation methodology: Normal, Special, and Agent. "Normal" evaluates tool usage in basic scenarios; "Special" evaluates tool usage in situations with ambiguous or incomplete instructions; "Agent" evaluates tool usage through multi-agent interactions to simulate real-world, multi-turn dialogues. We conducted extensive experiments using ACEBench, analyzing various LLMs in-depth and providing a more granular examination of error causes across different data types.

#### 1 Introduction

Large Language Models (LLMs), such as GPT-4 (Achiam et al., 2023), have demonstrated exceptional performance across numerous natural language processing tasks (Naveed et al., 2023; Qu et al., 2025; Mialon et al., 2023). Studies have shown that incorporating tools can significantly expand LLM capabilities, particularly in specialized domains such as mathematics (Das et al., 2024; Bulusu et al., 2024; Gou et al., 2023; Veerendranath

et al., 2024), programming (Xu et al., 2024), and reasoning (Chen et al., 2022; Shao et al., 2022; Surís et al., 2023; Yang et al., 2023). On one hand, integrating tools into LLMs can enhance capabilities in multiple domains, for example, ToolTransformer (Schick et al., 2023) enhances the ability of LLMs to solve complex problems by utilizing tools. On the other hand, adopting a tool usage paradigm can improve the robustness of the response and the transparency of the generation, thus increasing the explainability and trust of users (Schick et al., 2023), as well as improving the system's adaptability. As this field continues to evolve, it is essential to comprehensively evaluate all aspects of tool usage, particularly in complex scenarios.

While several studies have focused on evaluating tool usage (Yan et al., 2024; Guo et al., 2024; Wang et al., 2024a; Qin et al., 2023; Wang et al., 2024b; Zhuang et al., 2023; Lu et al., 2024), there are still some shortcomings in the existing tool-use benchmarks. Firstly, existing benchmarks lack multi-turn dialogue evaluation in real-world scenarios. For example, the multi-turn dialogues in BFCL (Yan et al., 2024) and HammerBench (Wang et al., 2024a) are composed of predefined fixed content combinations. Secondly, current tool-use benchmarks (Qin et al., 2023; Guo et al., 2024; Huang et al., 2023; Li et al., 2023) lack fine-grained evaluation and personalized data assessment. Additionally, existing benchmarks (Qin et al., 2023; Guo et al., 2024; Wang et al., 2024b) ignore the assessment of special cases, or the evaluation methods are simplistic (Yan et al., 2024), as user instructions in real life are not always perfect(Wang et al., 2024c). The model's ability to recognize and handle these issues is also crucial for evaluation. Lastly, evaluation costs are high (Qin et al., 2023; Guo et al., 2024), as many studies rely on advanced large models for evaluation.

To address these shortcomings, we propose

<sup>&</sup>lt;sup>†</sup>Equal Contributions. Work was done during an internship at Huawei Noah's Ark Lab. \*Corresponding authors.

<sup>\*</sup>The code are already publicly available at GitHub.

Table 1: Comparison of benchmarks across different evaluation criteria. "LLM-Free" refers to result evaluation without relying on LLMs. "Robustness" refers to incomplete or unclear user instructions. "Interactiveness" refers to the dynamic interaction between the model and the environment. "Atomic-Level" refers to analyzing from the atomic-level capabilities. "Personalization" refers to the inclusion of personal likes.

Benchmark	LLM-Free	Robustness	Interactiveness	Atomic-Level	Personalization
MetaTool (Huang et al., 2023)	✓	Х	Х	Х	Х
API-Bank (Li et al., 2023)	✓	X	×	×	×
Stable ToolBench (Guo et al., 2024)	X	X	×	×	×
BFCL (Yan et al., 2024)	✓	✓	×	×	×
$\tau$ -Bench (Yao et al., 2024)	✓	X	✓	×	×
HammerBench (Wang et al., 2024a)	X	✓	×	×	×
ACEBench (Ours)	✓	✓	✓	✓	✓

ACEBench, a comprehensive tool-use benchmark that includes the following categories:

**Normal.** Consists of fixed question-answer pairs and encompasses a variety of scenarios, including single-turn dialogues, multi-turn dialogues, and personalized scenario data. It also includes evaluations of atomic-level capabilities.

**Special.** Includes imperfect instructions, such as instructions containing incomplete parameters, incorrectly formatted parameters, or questions irrelevant to the capabilities of the candidate functions.

**Agent.** Encompasses real-world scenarios, abstracted to construct multi-turn, multi-step tool invocation scenarios, divided into multi-turn and multi-step cases depending on whether the user participates in the dialogue process.

The three categories above cover most of the tool usage scenarios for LLMs, and detailed explanations of each category can be found in Appendix A. Our main contributions are as follows:

- Comprehensive Benchmark Evaluation. We propose a comprehensive benchmark for evaluating LLMs' tool usage, covering various scenarios, including more fine-grained evaluation perspectives and assessments under imperfect instructions and providing more stable evaluation metrics.
- Sandbox Environment and Automated Evaluation System. We build an end-to-end automated evaluation system and develop a sandbox environment construction scheme for multi-turn, multi-step tool invocation based on real-world scenario abstraction.
- Extensive Experimental Validation. Through extensive experiments, we demonstrate that our benchmark provides a more comprehensive analysis with greater distinction, offering a clearer

evaluation of LLMs' tool usage.

#### 2 Related Works

The emerging trend of leveraging LLMs' tool-use capabilities in real-world applications underscores the need for comprehensive evaluations of their performance and effectiveness. Despite recent advancements, existing benchmarks for evaluating the tool-use capabilities of LLMs still have significant limitations

Stable ToolBench (Guo et al., 2024) addresses the issue of unstable external APIs by employing a virtual API server, but its dependence on large models for evaluation results in high costs and scalability challenges. BFCL (Yan et al., 2024) introduces a benchmark for tool use in multi-turn dialogue scenarios. Yet, it assembles dialogues from fixed content, failing to capture the dynamic and adaptive nature of real-world interactions. Similarly,  $\tau$ -Bench (Yao et al., 2024) evaluates language agents' ability to engage with human users while adhering to domain-specific rules. Still, its narrow focus on just two scenarios limits its generalizability across diverse tasks. HammerBench (Wang et al., 2024a) improves upon this by incorporating datasets derived from popular mobile applications and merging dialogues to simulate typical question-answer trajectories. However, like BFCL, its multi-turn dialogues are simplistic concatenations of pre-defined content, which do not reflect the complexities of real-world conversational dynamics. In addition, some benchmarks (Qin et al., 2023; Guo et al., 2024) rely on large language models (LLMs) for result evaluation, leading to high costs and unstable operations.

In contrast, our work addresses these limitations by expanding the scope of evaluation to encompass a broader range of tool usage scenarios. We

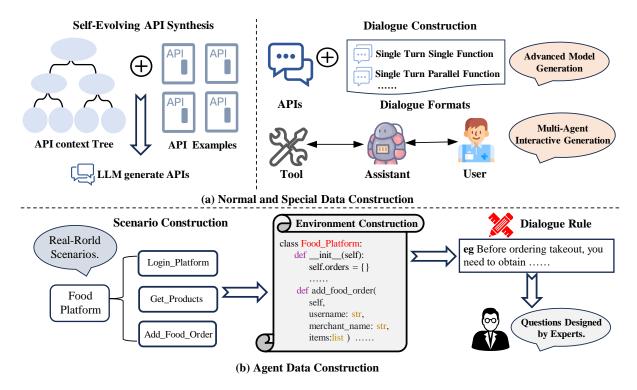


Figure 1: Dataset construction pipeline. (a) Normal and Special data construction: API synthesis module (left), Dialogue generation module (right). (b) Agent Data Construction: include scenario construction, environment construction and question design.

propose a framework that simulates realistic multiturn dialogue processes and enables end-to-end automated assessment, thereby reducing evaluation costs and improving scalability. A comparative analysis of ACEBench against recent benchmarks, as shown in Table 1, demonstrates its effectiveness in overcoming these challenges.

#### 3 ACEBench

#### 3.1 Dataset

We constructed two linguistically parallel versions of the dataset (Chinese and English), ensuring equal distribution of data types between them. The final dataset comprises 2,000 annotated entries.

#### 3.1.1 Data Construction

The Normal and Special data are automatically generated by LLMs, whereas the Agent data is constructed by experts. Creation details for some data are provided in Appendix Section B.

Normal and Special Data Construction. We employ a fully automated LLM-based generation pipeline specifically designed for Normal and Special Data, as illustrated in Figure 1(a). The API context tree is a hierarchical structure that organizes APIs from broad domains into increasingly

fine-grained functional categories. For example, under the domain of *Education*, we define subdomains such as *Subject Teaching* and *Learning Analytics*, where the branch *Subject Teaching* further includes APIs like generate\_math\_quiz or evaluate\_student\_progress.

- (1) API Synthesis. We use real APIs from various real-world scenarios as reference during construction to enhance authenticity. To ensure the stability of the data, we use synthetic APIs to construct the evaluation dataset, referencing real-world APIs as a guide. We employ a self-evolution approach by building a hierarchical API context tree to ensure the generated APIs cover a wide range of domains and functionalities (Liu et al., 2024b). Initially, we extract relevant information from technical documents to guide the API generation. As the process progresses, the context tree is gradually expanded, ultimately ensuring the depth and breadth of the generated APIs.
- (2) Dialogue Construction. We use two different dialogue generation pipelines built on the constructed API pool from which three to six candidate APIs are selected for each evaluation instance. For most cases, APIs are chosen randomly. However, for instances requiring specific functionality (e.g., similar APIs or multi-turn scenarios), advanced

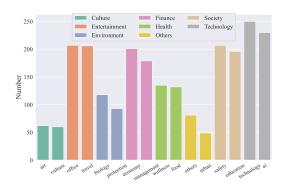


Figure 2: Distribution of APIs in terms of domains (Top 2 subcategories for each category).

methods, including graph-based sampling (Wang et al., 2024d), are used. Simple cases or those with predefined functionality use a template-based generation, where a single generator produces dialogues to ensure consistency. We employ a multiagent dialogue pipeline for more complex scenarios, where three agents (user, assistant, and tool) to simulate real-world interactions. Both pipelines are supported by carefully hand-crafted examples to ensure comprehensive coverage and diversity. Importantly, our benchmark explicitly supports function composition, both within a single turn and across multiple turns. For instance, a system may invoke multiple APIs simultaneously in one turn or sequentially over multiple turns, depending on the task complexity.

**Agent Data Construction.** We implement a carefully curated human-expert construction framework specifically tailored for Agent Data generation, as shown in Figure 1(b).

- (1) Scenario Construction. Through systematic abstraction of real-world interaction scenarios (such as food delivery services and telecommunication operations), we design functional modules with well-defined business semantics and specify each scenario's core state variables (e.g., order status, account balance) and intrinsic property sets.
- (2) Sandbox Environment Construction. We constructed an isolated sandbox environment with three core components: (1) task-specific scenarios (e.g., food ordering, travel booking, telecommunication operations) with clearly defined objectives and interaction flows; (2) executable functions aligned with these scenarios, each with explicit input/output specifications and preconditions; and (3) a user simulator powered by large language models, which engages in multi-turn dialogues and provides dynamic feedback based on system ac-

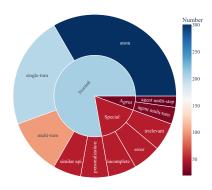


Figure 3: Visualization of the data composition of ACEBench.

tions (e.g., "Insufficient balance" or "I would like to add a drink").

(3) Question Design. Based on predefined multiturn dialog specifications tailored to different scenarios, domain experts systematically crafted the conversational questions through an iterative annotation process.

#### 3.1.2 Multi-Stage Data Verification

To address issues like mismatched answers or ambiguous criteria, we have implemented a multistage verification process.

Automated Quality Inspection. The data first undergoes a rule-based quality inspection module, which evaluates four dimensions: clarity of API definitions, executability of function calls, accuracy of dialogues, and consistency of data samples, effectively filtering out formatting and spelling errors. Next, the data enters the model-based quality verification module, which uses LLMs to detect semantic errors, employing a voting mechanism to ensure consistency in evaluation.

Human Quality Inspection. In the initial evaluation, the dataset remaining after automated quality inspection is assessed by three LLMs to assist human experts in data screening. Valid data is retained, while potentially problematic data is placed in the error candidate pool. These flagged entries undergo a two-step expert review process, where two experts independently assess and suggest modifications, and a third expert consolidates feedback, revising problem statements, API definitions, and answers. The revised data is re-evaluated and manually verified, and three rounds of optimization are performed to ensure a high-quality dataset.

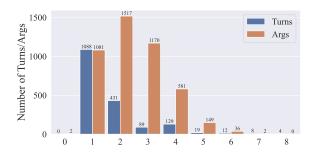


Figure 4: Distribution of dialogue turns and API argument numbers.

#### 3.1.3 Data Analysis

To demonstrate the breadth and comprehensiveness of ACEBench, we provide a detailed analysis of its test case distributions. Specific examples of each data type can be found in Appendix C.

**Domain of APIs.** The ACEBench API boasts a comprehensive coverage of 8 major domains and 68 sub-domains, spanning various aspects of daily life, including technology, finance, entertainment, society, health, culture, environment, and others. It offers a rich collection of 4,538 APIs in both Chinese and English. The distribution of these APIs is visualized in the accompanying Figure 2.

**Data Composition.** ACEBench consists of three categories of test samples: Normal, Agent, and Special, where each category is divided into several subcategories. The data composition is visualized in Figure 3, demonstrating a comprehensive coverage of tool-use capabilities, from simple single-turn tool invocations to complex multi-turn interactions involving users and environments. They include scenarios requiring multiple steps and interactions with the environment, as well as cases where tool calls are infeasible.

Number of turns and arguments. The test data in ACEBench covers a wide range of complexities. Specifically, we statistically analyzed the number of dialogue turns and the number of arguments in the called apis, which are visualized in Figure 4. The results show that the number of dialogue turns ranges from 1 to 8, encompassing most real-world scenarios. These samples with varying numbers of turns and arguments further form a test suite that covers a broader range of difficulties.

#### **3.2** Eval

In this section, we introduce the evaluation process, and the inference prompt for evaluation can be found in Appendix D.

#### 3.2.1 Normal Evaluation

As shown in the left part of Figure 5, we evaluate Normal Data by comparing the model's function call output with the ground truth using AST parsing. For cases with multiple valid answers, we employ a candidate answer pool where matching any candidate constitutes correctness. Evaluation uses Accuracy metric (1=full match, 0=mismatch).

#### 3.2.2 Special Evaluation

As illustrated in the middle section of Figure 5, the evaluation of Special Data primarily assesses the model's capability in problem identification. Specifically, the model must: (1) detect and alert missing parameters, (2) accurately locate erroneous parameters, and (3) recognize task-function mismatches. For each case, Accuracy is scored as 1 if correctly identified, otherwise 0.

#### 3.2.3 Agent Evaluation

As shown in the right part of Figure 5, we evaluate the agent's capabilities by assessing the model's proficiency in utilizing tools during human-agent interactions, employing gpt-40 as a user simulator for testing purposes. There are two evaluation metrics:

**End-to-End Accuracy** is evaluated by comparing the instance attributes of the corresponding class with the target. If all attributes match exactly, the Accuracy is 1; otherwise, the Accuracy is 0.

**Process Accuracy** is determined by the consistency between the actual function call process and the ideal process. It is expressed as  $\frac{n}{m}$ , where m represents the ideal function call process, and n represents the degree of match between the actual and ideal processes.

#### 3.2.4 Overall Accuracy

The Overall Accuracy is computed as a weighted sum of the accuracies for the Normal, Special, and Agent data types, where the weights are determined by the square roots of their respective sample sizes. The details can be found in Appendix E.

#### 4 Experiments

In this section, we present a comprehensive set of experiments designed to evaluate the performance of LLMs on ACEBench.

**Experimental Setup.** In our evaluation, we examine seven closed-source LLMs, including the GPT-4 series (Achiam et al., 2023), Qwen-Max (Yang et al., 2024), Gemini-1.5-Pro (Team et al., 2024),

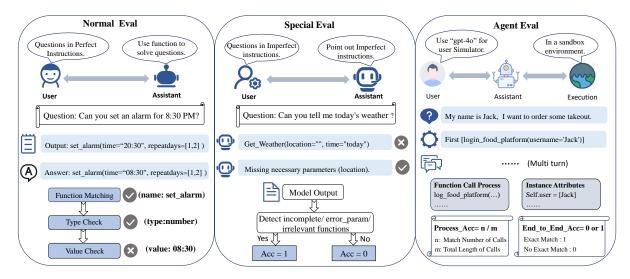


Figure 5: Overview of evaluation process: The left represents 'Normal' evaluation: AST-based function and parameter verification. The middle illustrates 'Special' evaluation: Imperfect instruction defect diagnosis. The right shows 'Agent' evaluation: State transition analysis via user-model interaction

Claude-3.5-Sonnet (Anthropic, 2024), and Doubao-Pro-32K (ByteDance, 2025). Additionally, a wide range of open-source language models are assessed, such as the Qwen2.5 series (Yang et al., 2024), Llama3 series (Dubey et al., 2024), Phi-3-Mini (Abdin et al., 2024), Deepseek-V3(Liu et al., 2024a), and DeepSeek-Coder-V2 (Zhu et al., 2024). Furthermore, four tool-learning-enhanced models were evaluated: Hammer2.1-3B, Hammer2.1-7B (Lin et al., 2024), xLAM-7B-r (Liu et al., 2024c), and Watt-Tool-8B (Watt-AI, 2024).

#### 4.1 Main results and analysis

The comprehensive experimental results for the Chinese and English datasets are presented in Table 2, with detailed results for each language provided in Appendix F. We can draw the following important conclusions:

#### General Conclusion on Model Performance.

The overall best performance remains dominated by closed-source models, such as the GPT-4 series. However, the performance gap between certain open-source models, such as Qwen2.5-Coder-32B-Instruct, Qwen2.5-72B-Instruct and DeepSeek-V3, and their closed-source counterparts is progressively narrowing. This trend suggests that open-source models are steadily catching up to closed-source models, driven by advancements in model architecture and training methodologies.

Loss of Generalization in Fine-Tuned Models. As shown in Table 3, models fine-tuned on specific datasets, such as Watt-Tool-8B (Watt-AI, 2024),

xLAM-7B (Liu et al., 2024c), and Hammer2.1-7B(Lin et al., 2024), exhibit a significant decline in performance on the Special dataset. This decline can primarily be attributed to the fact that while fine-tuning enhances a model's performance on specialized tasks, it can also lead to a loss of generalization, making the model less effective on new or broader instruction-following tasks. This phenomenon highlights the importance of balancing task-specific performance and generalization capability during model optimization.

Performance Limitations of Large Models in Complex Tasks. As shown in Table 4, most models exhibit an end accuracy of less than 50% on Agent data tasks. This can be attributed to the fact that completing such tasks in dynamic environments, which simulate real-world multi-turn interactions, requires more than just performing individual tool operations. The model must also integrate contextual information during tool usage and account for the interdependencies between tool calls, which significantly increases task complexity. Furthermore, these tasks demand advanced reasoning and adaptability, which even large models may struggle with due to the challenges of maintaining consistency across long-term interactions and responding to the evolving nature of the task.

#### 4.2 Error Analysis

**Error Analysis of Normal Data.** As shown in Figure 6, we observe from the error type distribution on Normal data that param value error dominate

Table 2: Comprehensive evaluation of different models on ACEBench for Chinese and English combined (%).

Model			N	ormal			Special	Agent	Overall
Model	Atom	Single-Turn	Multi-Turn	Similar API	Preference	Summary	Special		
		Closed	-Source Large	Language Mo	dels				
GPT-4o	93.4	84.5	77.0	85.0	83.0	87.6	93.0	63.8	85.4
GPT-4-Turbo	93.2	84.8	77.5	86.0	86.0	88.0	86.7	67.5	84.5
Qwen-Max	91.2	80.5	68.0	83.0	83.0	84.2	74.0	64.3	78.4
GPT-4o-Mini	86.5	76.0	66.5	77.0	78.0	79.9	79.0	33.3	72.5
Gemini-1.5-Pro	84.5	76.8	64.5	80.0	78.0	79.0	78.7	25.5	70.7
Claude-3-5-Sonnet	76.9	72.5	62.5	71.0	72.0	72.9	77.4	39.5	68.9
Doubao-Pro-32k	79.8	55.5	58.0	76.0	66.0	70.7	55.0	25.0	59.4
		Open-	Source Large	Language Mod	iels				
Qwen2.5-Coder-32B-Instruct	90.2	81.0	71.0	83.0	81.0	84.1	80.7	60.8	79.6
DeepSeek-V3	91.5	84.0	77.0	83.0	83.0	86.5	73.0	34.5	74.8
Qwen2.5-72B-Instruct	86.8	80.3	69.5	83.0	81.0	82.1	75.7	45.0	74.7
Llama-3.1-70B-Instruct	82.5	68.3	63.5	79.0	68.0	75.5	38.3	42.3	60.4
Qwen2.5-7B-Instruct	76.0	60.3	58.5	72.0	67.0	69.4	47.0	13.8	54.8
DeepSeek-Coder-V2-Lite-Instruct	75.2	57.8	46.5	72.0	65.0	66.4	40.3	2.0	49.5
Qwen2.5-Coder-7B-Instruct	76.0	63.8	57.5	74.0	68.0	70.1	22.3	15.5	48.9
Watt-Tool-8B	85.7	69.3	55.5	79.0	64.0	75.6	6.0	2.8	45.7
Hammer2.1-7B	73.7	57.5	40.0	62.0	55.0	62.8	14.7	16.8	42.9
Llama-3.1-8B-Instruct	51.9	39.8	28.0	66.0	46.0	46.6	21.0	5.3	33.4
Phi-3-Mini-128k-Instruct	57.2	39.3	23.0	58.0	32.0	46.5	18.7	0.8	32.0
xLAM-7B-r	43.5	22.0	19.0	61.0	0.0	33.7	2.7	8.8	21.6
Llama-3.2-3B-Instruct	38.7	15.3	9.0	42.0	32.0	29.6	9.4	0.0	19.6
Hammer2.1-3B	22.4	11.5	3.5	40.0	20.0	18.7	1.0	1.5	11.3

Table 3: The Accuracy of Models on Special Data

Model	Incomplete	Error	Irrelevant
Llama-3.1-8B-Instruct	29.0	20.0	14.0
Watt-Tool-8B	7.0	1.0	10.0
Qwen2.5-7B-Instruct	26.0	36.0	79.0
xLAM-7B-r	1.0	3.0	4.0
Llama-3.2-3B-Instruct	13.0	12.0	3.0
Hammer2.1-3B	0.0	3.0	0.0

across all models. This highlights the models' difficulty in generating specific values, likely due to limited contextual understanding and the complexity of numerical distributions. Output format error is the second most common, suggesting room for improvement in generating code that follows predefined formats and syntactic rules. These issues may stem from inconsistencies in training data and the models' limited ability to learn rule-based generation. In contrast, function name and param type errors are less frequent, indicating that the models excel in matching function calls and handling data types. While the models show strong function invocation abilities, further improvements are needed in numerical generation and format compliance. Specific error examples for Normal data can be found in Appendix G.1.

**Error Analysis of Special Data.** As shown in Table 5, we identified two main types of model errors:

Table 4: Performance evaluation of different models on Agent Data: **PA** represents Process Accuracy, **EA** represents End-to-End Accuracy (%).

Model	Multi	Turn	Multi Step		
NIOGO1	EA	PA	EA	PA	
GPT-4-Turbo	50.0	66.0	85.0	89.5	
DeepSeek-V3	31.5	54.5	37.5	53.0	
Claude-3-5-Sonnet	21.5	41.5	57.5	76.5	
DouBao-Pro-32k	20.0	45.5	30.0	47.5	
Qwen2.5-7B-Instruct	15.0	28.0	12.5	15.5	
Hammer2.1-7B	8.5	33.5	25.0	42.5	

The first type is "Error Detection", which refers to the model's complete failure to detect issues in the user's instructions or its inability to identify problems according to the prompt's formatting requirements. The second type is "Error Correction," where the model detects the problem but provides unclear feedback. For example, the model might indicate that there is an issue, but fails to specify which parameter values are incorrect or what critical information is missing. Results show that most errors in special-type scenarios are caused by "Error Detection", highlighting a critical gap in the model's problem-detection capabilities. This suggests that the model needs to learn not only simple tool invocation but also how to identify corresponding issues under imperfect instructions. Specific

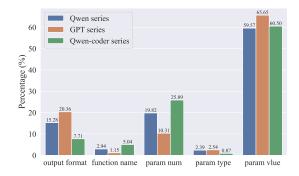


Figure 6: Error type distribution on Normal Data.

error examples can be found in Appendix G.2.

Error Analysis of Agent Data. Our analysis identifies three primary causes of Agent errors. First, function call errors occur when the model fails to select the appropriate function or provide parameters that do not meet the required specifications, reflecting a lack of understanding of tool-use capabilities and parameter constraints. Second, rule violations arise when the model disregards predefined scene rules, skipping necessary steps or breaking key task logic, highlighting deficiencies in its comprehension and execution. Finally, information mismanagement results from the model's inability to correctly record or process contextual information during multi-turn interactions, leading to outputs that diverge from expectations. As shown in Figure 21, we illustrate an error caused by missing information.

#### 4.3 Further Analysis

**Scaling Law.** We evaluated the performance of Qwen2.5-Coder (3B, 7B, 14B, 32B) and Qwen2.5-Instruct (3B, 7B, 14B, 32B, 72B) on the ACEBench dataset. As shown in Figure 7, the experimental results demonstrate that performance improves significantly across various tasks as the model size increases, with particularly strong results observed in high-complexity tasks. However,

Table 5: Error type distribution across different model series on Special Data.

Model	<b>Error Detection</b>	Error Correction
Watt-Tool-8B	188	4
Hammer2.1-7B	172	7
Phi-3-mini-128k-instruct	143	15
Qwen2.5-3B-Instruct	130	36
xLAM-7B-r	195	1
Llama-3.1-8B-Instruct	145	6
Hammer2.1-3B	197	0
Llama-3.2-3B-Instruct	166	9

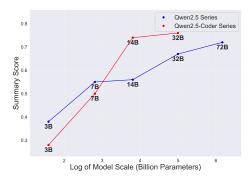


Figure 7: Scaling Law of LLMs on ACEBench.

it is worth noting that as the model size continues to grow, the rate of performance improvement begins to slow down, especially between the 32B and 72B models. This indicates that while increasing the model parameters brings substantial performance gains initially, the marginal benefits of scaling up further decrease, making additional improvements more challenging.

Table 6: Accuracy comparison of prompting strategies on English Normal Data (%).

Model	Standard	Condensed	Minimal
Qwen2.5-3B-Instruct	34.5	31.8	27.8
Qwen2.5-7B-Instruct	48.5	47.5	45.5
Qwen2.5-14B-Instruct	56.3	54.0	47.5

**Impact of Prompting Strategies.** Prompt design significantly affects language model performance. We tested three strategies (see Appendix D.4):

- (1)Standard Prompt: A comprehensive template designed to eliminate interference from information insufficiency, ensuring a fair evaluation.
- (2)Condensed Prompt: A compact version retaining core instructions, testing performance with reduced but sufficient guidance.
- (3)Minimal Prompt: A highly abbreviated form (e.g keywords) to assess the model's ability to infer tasks from ultra-concise input.

The experimental results in Table 4 demonstrate that models utilizing standard prompt templates achieve the highest overall accuracy. This optimal performance can be attributed to the rigorous formatting specifications in standard prompts, which effectively mitigate interference from extraneous variables. These empirical findings establish a positive correlation between prompt standardization and model performance, providing key insights for future prompt engineering: enhancing the standardization of function-calling prompts with explicit

formatting requirements can significantly improve execution accuracy.

#### 5 Conclusion

This paper introduces ACEBench, a comprehensive tool-use benchmark designed to evaluate the tool-use capabilities of Large Language Models (LLMs), including data from three types: normal, special, and agent. It addresses key limitations of existing evaluation benchmarks, such as the lack of multi-turn dialogue assessments in real-world scenarios, the absence of fine-grained evaluations for parameter-type function calls, and the high costs associated with using large models for evaluation. The experimental results indicate that models finetuned on specific tool-use datasets to struggle with generalization when faced with complex or imperfect instructions, and code capabilities enhance the tool-use performance of large models. Through extensive experiments, we demonstrate the effectiveness of ACEBench in providing deeper insights into the tool-use abilities of various models.

#### Limitations

We acknowledge several limitations in our evaluation of ACEBench for assessing the tool-use capabilities of large language models. Firstly, while our test data is generated by large language models and various measures have been taken to ensure its authenticity and diversity, a gap remains when compared to data from real-world applications. This discrepancy may impact the evaluation of the model's performance in real-world scenarios. Secondly, for the Agent data, the design of evaluation scenarios relies on manual construction, which somewhat limits the diversity and coverage of the evaluation framework. Finally, the current benchmark supports only structured JSON-style function calls. While this format has become the de facto standard in many mainstream LLM applications (e.g., OpenAI, Anthropic, Google), realworld tool invocation often involves other representations (e.g., Python code, XML, DSLs). Extending ACEBench to incorporate diverse invocation formats would provide a more comprehensive evaluation of tool-use capabilities, and we leave this as future work.

#### Acknowledgements

This work is supported by National Natural Science Foundation of China (62502310).

#### References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv* preprint arXiv:2404.14219.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.
- Anthropic. 2024. Claude 3.5 sonnet. Accessed: January 2025.
- Arya Bulusu, Brandon Man, Ashish Jagmohan, Aditya Vempaty, Jennifer Mari-Wyka, and Deepak Akkil. 2024. Mathviz-e: A case-study in domain-specialized tool-using agents. *arXiv preprint arXiv:2407.17544*.
- ByteDance. 2025. Doubao product page. Accessed: January 2025.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2022. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv* preprint *arXiv*:2211.12588.
- Debrup Das, Debopriyo Banerjee, Somak Aditya, and Ashish Kulkarni. 2024. Mathsensei: A toolaugmented large language model for mathematical reasoning. arXiv preprint arXiv:2402.17231.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv* preprint arXiv:2407.21783.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yujiu Yang, Minlie Huang, Nan Duan, Weizhu Chen, et al. 2023. Tora: A tool-integrated reasoning agent for mathematical problem solving. *arXiv preprint arXiv:2309.17452*.
- Zhicheng Guo, Sijie Cheng, Hao Wang, Shihao Liang, Yujia Qin, Peng Li, Zhiyuan Liu, Maosong Sun, and Yang Liu. 2024. Stabletoolbench: Towards stable large-scale benchmarking on tool learning of large language models. *arXiv preprint arXiv:2403.07714*.
- Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao Wan, Neil Zhenqiang Gong, et al. 2023. Metatool benchmark for large language models: Deciding whether to use tools and which to use. *arXiv preprint arXiv:2310.03128*.
- 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. arXiv preprint arXiv:2304.08244.
- Qiqiang Lin, Muning Wen, Qiuying Peng, Guanyu Nie, Junwei Liao, Jun Wang, Xiaoyun Mo, Jiamu Zhou, Cheng Cheng, Yin Zhao, et al. 2024. Hammer: Robust function-calling for on-device language models via function masking. *arXiv preprint arXiv:2410.04587*.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024a. Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437.
- Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, et al. 2024b. Toolace: Winning the points of llm function calling. *arXiv* preprint arXiv:2409.00920.
- Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, et al. 2024c. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets. *arXiv preprint arXiv:2406.18518*.
- Jiarui Lu, Thomas Holleis, Yizhe Zhang, Bernhard Aumayer, Feng Nan, Felix Bai, Shuang Ma, Shen Ma, Mengyu Li, Guoli Yin, et al. 2024. Toolsandbox: A stateful, conversational, interactive evaluation benchmark for llm tool use capabilities. *arXiv preprint arXiv:2408.04682*.
- Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, et al. 2023. Augmented language models: a survey. arXiv preprint arXiv:2302.07842.
- Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. 2023. A comprehensive overview of large language models. arXiv preprint arXiv:2307.06435.
- 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. arXiv preprint arXiv:2307.16789.
- Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2025. Tool learning with large language models: A survey. *Frontiers of Computer Science*, 19(8):198343.
- T. Schick, J. Dwivedi-Yu, R. Dessì, et al. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551.

- Zhihong Shao, Fei Huang, and Minlie Huang. 2022. Chaining simultaneous thoughts for numerical reasoning. *arXiv* preprint arXiv:2211.16482.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. In *Proceedings of the IEEE/CVF Interna*tional Conference on Computer Vision, pages 11888– 11898.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv* preprint arXiv:2403.05530.
- Vishruth Veerendranath, Vishwa Shah, and Kshitish Ghate. 2024. Calc-cmu at semeval-2024 task 7: Pre-calc-learning to use the calculator improves numeracy in language models. *arXiv preprint arXiv:2404.14355*.
- Jun Wang, Jiamu Zhou, Muning Wen, Xiaoyun Mo, Haoyu Zhang, Qiqiang Lin, Cheng Jin, Xihuai Wang, Weinan Zhang, and Qiuying Peng. 2024a. Hammerbench: Fine-grained function-calling evaluation in real mobile device scenarios. arXiv preprint arXiv:2412.16516.
- Pei Wang, Yanan Wu, Zekun Wang, Jiaheng Liu, Xiaoshuai Song, Zhongyuan Peng, Ken Deng, Chenchen Zhang, Jiakai Wang, Junran Peng, et al. 2024b. Mtu-bench: A multi-granularity tool-use benchmark for large language models. *arXiv* preprint arXiv:2410.11710.
- Wenxuan Wang, Juluan Shi, Chaozheng Wang, Cheryl Lee, Youliang Yuan, Jen-tse Huang, and Michael R Lyu. 2024c. Learning to ask: When Ilms meet unclear instruction. *arXiv preprint arXiv:2409.00557*.
- Z. Wang, X. Zeng, W. Liu, and et al. 2024d. Toolflow: Boosting Ilm tool-calling through natural and coherent dialogue synthesis. *arXiv* preprint *arXiv*:2410.18447.
- Watt-AI. 2024. Watt-tool 8b. Accessed: January 2024.
- Shuyuan Xu, Zelong Li, Kai Mei, and Yongfeng Zhang. 2024. Core: Llm as interpreter for natural language programming, pseudo-code programming, and flow programming of ai agents. *arXiv preprint arXiv:2405.06907*.
- Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. 2024. Berkeley function calling leaderboard. https://gorilla.cs.berkeley.edu/blogs/8\_berkeley\_function\_calling\_leaderboard.html. Accessed: 2024-01-16.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.

- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023. Mmreact: Prompting chatgpt for multimodal reasoning and action. *arXiv* preprint arXiv:2303.11381.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024.  $\tau$ -bench: A benchmark for toolagent-user interaction in real-world domains. *arXiv* preprint arXiv:2406.12045.
- Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. 2024. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv preprint arXiv:2406.11931*.
- Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. Toolqa: A dataset for llm question answering with external tools. *Advances in Neural Information Processing Systems*, 36:50117–50143.

#### **A** Detailed Descriptions of Test Cases

#### A.1 Data Categories Description

We divide the benchmark into three main categories: Normal, Special, and Agent. Below is a detailed description of each category.

#### **Normal Data**

The Normal Data consists of fixed question-answer pairs, where each question corresponds to a correct function call. It is categorized into the following categories: Single-Turn, Multi-Turn, Similar APIs, Preference, and Atom.

**Single-Turn**: There is only one round interaction between the user and the assistant, and based on the number of function calls in the response, it is divided into single-turn single function calls and single-turn parallel function calls.

**Multi-Turn**: There are multiple interactions between the user and the assistant. The conversation can be categorized into two types: (**switch**) The conversation progresses by changing topics. (**adjust**) The conversation evolves by refining or modifying the original question.

**Similar APIs**: The candidate APIs exhibit significant similarity, particularly focusing on the same topic. This similarity presents a challenge for the assistant, requiring it to effectively distinguish between the APIs and accurately select the most appropriate ones.

**Preference**: Besides the candidate APIs, the assistant is provided with supplementary user profile data. This type of information necessitates the assistant's ability to mine user-specific factors, such as past interactions, interests, or other personalized attributes, to generate argument values.

**Atom**: Atom Data refers to a set of APIs that contain only specific parameter types, such as candidate functions where the parameters exclusively involve numbers, lists, etc. This design is intended to explore whether the type of function parameters affects the model's ability to handle data filling. We have divided the Atom data into five types: number, enum, list, bool, and object.

#### **Special Data**

The Special Data refers to situations where the model is unable to resolve the problem posed in the instruction using the candidate functions(Wang et al., 2024c). It is categorized into the following categories: Incomplete, Error, and Irrelevant.

**Incomplete**: Refers to situations where the key information required for the function call is missing in the query, such as the absence of "required" parameters.

**Error**: Refers to situations where the instruction contains parameters or names that do not meet the required format or constraints, such as matching a specific pattern or being selected from a predefined list, causing the function call to fail.

**Irrelevant**: Refers to situations where the instruction exceeds the function's capabilities, meaning none of the candidate functions can resolve the issue.

#### **Agent Data**

Agent Data refers to scenarios where completing a task in an environment modeled after real-world situations typically requires multi-step collaboration. In this study, we employ the GPT-40 language model to simulate user roles and replicate real-world interaction processes, thereby evaluating the model's performance in complex interactive settings. The key scenarios are defined as follows:

**Multi-step Scenario**: The user participates in only a single interaction throughout the entire dialogue flow.

**Multi-turn Scenario**: The user engages in multiple interactions across the entire dialogue cycle.

Agent Data currently encompasses the following fundamental scenarios:

(1)Mobile Application Simulation. The mobile application scenario provides digital lifestyle functionalities including communication services, integrated reminder and memo management systems, and alarm configuration capabilities. This environment simulates core smartphone operations with particular attention to notification handling and scheduling precision.

(2)Food Delivery Platform. This scenario simulates the core functionalities of a food delivery platform, primarily including merchant search, product browsing and ordering, order status tracking, and cancellation processing. The system implements essential operational procedures from merchant selection to order completion, supporting users throughout the entire food ordering experience.

(3)Financial Services Scenario. This module provides fundamental banking service simulations, primarily including: deposit/withdrawal transactions, account balance inquiries, fund transfers, and other routine banking operations, while also

supporting loan applications and repayment processes. The system maintains detailed transaction records and can generate basic financial statements, replicating the core services of real banking systems.

(4)Travel Booking Platform. This scenario simulates a standard ticketing system, enabling users to complete end-to-end operations including flight/train ticket inquiry, booking, payment, rescheduling, and cancellation. The system incorporates fare checking, seat selection, and order management functionalities, capable of handling itinerary changes and related ticket adjustments. It covers the complete user journey from search to ticket issuance.

Our team is actively working on designing more functional scenarios to enhance the platformin Agent Data.

#### B The Construction of Data

#### **B.1** Agent Data

As shown in Figure 8, the construction of Agent data can be summarized in the following steps:

First, through an in-depth analysis of real-world scenarios, extract key task requirements and modularize them into sub-scenarios in different domains (such as flight booking, food delivery platforms, and financial services), clearly define the specific functional objectives of each module.

Second, implement logical abstraction for each function through code, designing core processes such as user authentication, cost calculation, payment processing, and information recording. This ensures the code logic's scalability and robustness while comprehensively addressing exception handling.

Third, design interaction rules based on specific scenario requirements to standardize the interaction process between users and the Agent, such as verifying account, merchant information, and balance status in food delivery orders, and providing user guidance in exceptional cases.

Finally, combine real-world demands to design question formats and solutions, enabling the model to accurately meet user task requirements and achieve the desired outcomes.

#### **B.2** Special Data

#### **Irrelevant**

The irrelevant data refers to situations where the instruction exceeds the function's capabilities. And the construction method for Irrelevant data is similar to that of Normal data, where we simply remove the correct API from the candidate APIs of the common data.

#### **Incomplete**

The special data with incomplete instructions primarily refers to situations where key information is missing from the user's instructions, causing the function to be called incorrectly. Our main prompt for constructing incomplete data is shown in Figure 9. From the generated dialogue, we can extract the corresponding data and reference answers (the missing data). Next, we can convert the conversation we obtained into data. A specific example is shown in Figure 10.

```
Please refer to the example and continue the dialogue
based on the given tool definition and the beginning of
the conversation. The requirements are as follows:
1 The user's request is meaningful, requiring the use of
one tool, and the tool will be called once.
2 Calling the tool to fulfill the user's request still lacks
one or more required parameters.
Here is an example for reference
<tool_definition>
.
"name": "calculate_triangle_area",
"description": "Given the base and height of a
               triangle, calculate its area.",
"required": ["base", "height"]
<dialogue_example>
[User]: I want to calculate the area of a
triangle with a base of 5 and a height of 10.
[Assistant]:
[Thought] Missing parameters: base|height
[Response]:
Please provide both the base and the
height of the triangle.
[User]: The base is 5, and the height is 10.
[Assistant]:
<tool_usage>calculate_triangle_area
|{"base": 5, "height": 10}</tool_usage>
Here is the dialogue continuation you need to write:
<Tool Definition>
{tool_definition}
.....
```

Figure 9: Main prompt of Incomplete Data construction.

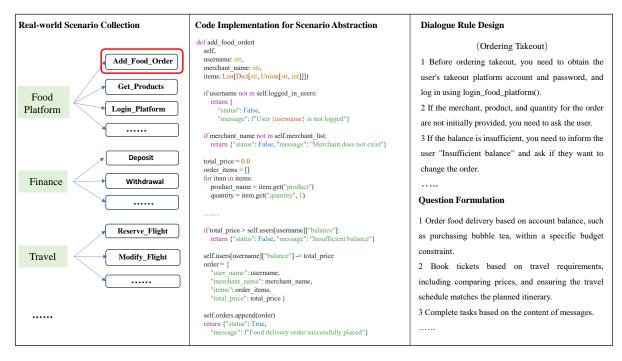


Figure 8: The construction of Agent Data. The left shows scenario sampling from real-world cases, the middle demonstrates the implementation of code tailored to specific scenarios, and the right presents examples of dialogue rules and question design for the scenarios.

```
candidate function
{
  "name": "book_flight",
  "description": "Flight booking ...",
  "arguments": {},
  "required": ["from_city_name"
   "to_city_name", "depart_date"]
}
obtained dialouge
<user> I plan to travel to Beijing.
Help me book a flight to Beijing.
</user> [Thought] Missing parameters:
from_city_name
[Response] May I know the departure
city you want to book?
<user> Shenzhen
converted data
[question] I plan to travel to Beijing.
Help me book a flight to Beijing.
[answer] Miss necessary parameter
(from_city_name) from (book_flight)
```

Figure 10: Specific example of Incomplete Data construction.

```
Please refer to the example and continue the dialogue
based on the given tool definition and the beginning of
the conversation. The requirements are as follows:
1 The user's request is specific and meaningful, requir-
ing one tool, and the tool should be called once.
2 The user's first question contains incorrect parameters
(i.e., it does not conform to the function definition's
pattern or format requirements).
Here is an example for reference
<tool_definition>
{"name": "TennisTeamDetails";
"description": "Retrieve detailed
information about a tennis team .",
"parameters": {
   "properties": {
    "team_name": {.....
"pattern": "^[a-zA-Z\\s]+$"}}}
<dialogue_example>
<user> Can vou retrieve the details of
     the tennis team named 'Team@1234'?
</user> [Thought] Error Parameter: 'Team@1234'
[Response] The team name you provided does
          not meet the naming criteria.
<user> Oh, I made a mistake.
   The team name is TeamABC.
</user> <tool_usage>TennisTeamDetails
|{"team_name": "TeamABC"}</tool_usage>"
Here is the dialogue continuation you need to write:
<Tool Definition>
{tool_definition}
```

Figure 11: Main prompt for Error Data construction.

#### Error

The special data which has error instructions mainly refers to situations where the instruction contains parameters or names that do not meet the required format or constraints, and the construction is shown in 11. Next, we can convert the conversation we obtained into data. A specific example is shown in 12.

```
candidate function
{ "name": "FootballTeamDetails",
  "description": "Retrieve information
  about a football team by its name.",
  "parameters": {
      "team_name": {.....
        "pattern": "^[a-zA-Z\\s]+$"}}}
obtained dialouge
<user> I want to know information about
the football team football$156.
</user> [Thought] Error Parameter:
         football$156.
[Response] The team name you provided
does not meet the naming criteria.
<user> .....
converted data
[question] I want to know information
about the football team footbaoo$156.
[answer] There is incorrect value
(football$156) for the (team_name).
```

Figure 12: Specific example of Error Data construction.

#### C Examples of Dataset

#### **C.1** Normal Examples

**Single-Turn.** The example of Normal Single-Turn Data is shown in Figure 13.

**Multi-Turn.** The example of Normal Multi-Turn Data is shown in Figure 14.

**Preference.** The example of Normal Preference Data is shown in Figure 15.

**Similar APIs.** The example of Normal Similar APIs Data is shown in Figure 16.

**Atom.** The example of Normal Atom Data is shown in Figure 17.

#### C.2 Special Examples

**Incomplete.** The example of Special Incomplete Data is shown in Figure 18.

**Error.** The example of Special Error Data is shown in Figure 19.

**Irrelevant.** The example of Special Irrelevant Data is shown in Figure 20.

#### C.3 Agent Examples

The example of Agent Data is shown in Figure 21 and Figure 22.

#### **D** Evaluation Inference Prompts

#### **D.1** Normal Prompt

The main evaluation inference prompt for Normal Data is shown in Figure 23 and Figure 24.

#### **D.2** Special Prompt

The main evaluation inference prompt for Special Data is shown in Figure 27.

#### **D.3** Agent Prompt

An example of the evaluation inference prompt for Agent Data in a specific scenario is shown in Figure 28. And an inference prompt is shown in Figure 29.

#### **D.4** Different Prompt Strategies

Standard Prompt is shown in Figure 23. Condensed Prompt is shown in Figure 25. Minimal Prompt is shown in Figure 26.

#### E Formula for Overall Accuracy

The formula for calculating the **Overall Accuracy** can be expressed as:

All 
$$Acc = A \cdot Acc_{Normal} + B \cdot Acc_{Special} + C \cdot Acc_{Agent}$$

where the coefficients A, B, and C are defined as:

$$A = \frac{\sqrt{n_{\text{Normal}}}}{\sqrt{n_{\text{Normal}}} + \sqrt{n_{\text{Special}}} + \sqrt{n_{\text{Agent}}}}$$

$$B = \frac{\sqrt{n_{\rm Special}}}{\sqrt{n_{\rm Normal}} + \sqrt{n_{\rm Special}} + \sqrt{n_{\rm Agent}}}$$

$$C = \frac{\sqrt{n_{\rm Agent}}}{\sqrt{n_{\rm Normal}} + \sqrt{n_{\rm Special}} + \sqrt{n_{\rm Agent}}}$$

where  $n_{\text{Normal}}$ ,  $n_{\text{Special}}$ ,  $n_{\text{Agent}}$  are the sample sizes for the **Normal**, **Special**, and **Agent** data types, and Accuracy refers to the accuracy of each respective category.

#### F Detailed Results

The accuracy evaluation of different models on English Data is shown in Table 7, and the accuracy evaluation of different models on Chinese Data is shown in Table 8.

#### **G** Error Examples

#### **G.1** Error Examples of Normal Data

**Wrong Function Name.** An example of wrong function name is shown in Figure 30.

**Wrong Param Type.** An example of wrong param type is shown in Figure 31.

**Wrong Param Value.** An example of wrong param value is shown in Figure 32.

**Wrong Output Format.** An example of wrong output format is shown in Figure 33.

#### **G.2** Error Examples of Special Data

**Error Detection** Figure 34 shows an example of error detection.

**Error Correction** Figure 35 shows an example of error correction.

### (question) Please assess the system with asset ID 12345, vulnerability ID 67890, with a severity level of high, and a potential threat that may lead to data leakage.

## (candidate function) { "name": "Threat\_Vulnerability\_Assessment\_Tool", "description": "Assess potential security threats and provide mitigation suggestions", "parameters": ...... } { "name": "userPrivacySettingsManager", "description": "Manage user's privacy", "parameters": ......

Figure 13: An example of Normal Single-Turn Data.

```
(conversation)
User: I want to analyze the pest and disease risks faced by my rice crop, which is currently in the tillering
system: Please provide the current temperature and humidity conditions user: The temperature is 28°C, and the humidity
is 80%.
tool:
[PestInfestationRiskAnalyzer(
    crop_information={"crop_type": "rice", "growth_stage": "tillering"},
    climate_conditions={"temperature": 28, "humidity": 80})]
system: Based on the analysis, your rice is facing a moderate risk of pests and diseases under the current climatic
conditions. Please take necessary preventive measures.
User: I also want to know the maintenance requirements for my production equipment. The device ID is Device 123.
system: Please provide the sensor readings of the device, including temperature, vibration, and humidity.
User: The device currently reads a temperature of 75°C, vibration level of 5, and surrounding humidity of 50%.
(candidate function)
  "name": "PestInfestationRiskAnalyzer",
  "description": "Analyze the pest and disease risk faced by the farmland ",
  "parameters": .....
}
{
  "name": "Predictive_Maintenance_AI_Tool",
  "description": "Predicting the maintenance needs of equipment using sensor data",
  "parameters": .....
  "name": "PersonalityTraitAnalyzer",
  "description": "Analyze an individual's personality traits",
  "parameters": .....
}
```

Figure 14: An example of Normal Multi-Turn Data.

```
(question)
Can you update my preferred notification method and also check if my current email address and home location are
properly updated in the system?
(candidate function)
  "name": "updateOrderStatusAlerts",
  "description": "Sends automated alerts to users regarding the status of their current orders.",
  "parameters": .....
{
  "name": "submitProductReview",
  "description": "Allows users to submit a review for a product they have purchased.",
  "parameters": .....
{
  "name": "updateUserProfile",
  "description": "Updates the user's profile information based on provided data.",
  "parameters": .....
}
(profile)
{
  "basic_features": {
    "UserName": "Michael Smith",
    "UserEmail": "mike.smith@example.com",
    "UserHomeLocation": "Los Angeles, CA",
    "UserBirthday": "1978-04-23",
"UserLanguage": "Spanish",
"UserTimeZone": "PST",
  "user_history": {
    "shopping": [
      "Searched for 'Nike running shoes' on app",
      "Added Nike Air Max to cart",
      "Checked coupon availability for Nike products",
      "Filtered search by 'Outdoor Equipment' category"
      "Selected 'High spending' filter for items over $500",
    "takeout": [
      "Ordered Chicken Fajitas on the takeout app for lunch",
      "Opted to receive promotional deals via phone calls",
      "Chose Debit Card ending in 5678 for payment on the takeout app",
    ]
 }
}
```

Figure 15: An example of Normal Preference Data.

# (question) My baby has had a visible vein on her nose for 5 days, and she's been crying a lot with a decreased appetite. Can you help? (candidate function) { "name": "baby\_health\_check\_A", "description": "Checks the common reasons for baby's persistent vein visibility on the nose and suggests actions. This API considers factors like skin thinness, crying, or overexertion", "parameters": ...... } { "name": "baby\_health\_check\_B", "description": "Examines baby's vein visibility and recommends seeing a doctor. Focuses on persistent visibility and associated symptoms", "parameters": ...... }

Figure 16: An example of Similar API Data.

```
(question)
I need a design for my new website. It's for a technology company focusing on user engagement.
(candidate function)
  "name": "WebDesignAssistant_generateDesign",
  "description": "Generates a website design based on industry and user experience focus.",
  "parameters": {
    "type": "object",
    "properties": {
      "industry": {
         "description": "The industry for which the website is being designed.",
         "type": "string",
         "enum": [
           "Technology"
          "Healthcare",
          "Education",
          "Finance"
        ],
        "default": "Technology"
      },
"userExperience": {.....}
    "required": ["industry", "userExperience"]
  }
}
{....}
```

Figure 17: An example of Atom (enum) Data.

```
(question)
I'm considering relocating my business to the Middle East. Can you provide me with a list of major cities?

(candidate function)
{"name": "Get_Middle_East_Cities",
    "description": "Retrieves a list of cities in the Middle East, sorted by overall score by default.",
    "parameters": {
        "properties": {
            "description": "The sorting order for the list of cities.",
            "type": "string",
            "enum": ["asc", "desc"],
        }
      },
      "required": ["sort"]
}}
```

Figure 18: An example of Special Incomplete Data.

Figure 19: An example of Special Error Param Data.

```
(question)

Could you help me find available restaurants in New York City?

(candidate function)
{
    "name": "Get_Weather_Report",
    "description": "Retrieve the current weather report for a specified location",
    "parameters": ......
}
{
    "name": "GetTravelDestinationInfo",
    "description": "Retrieves information about a specific travel destination",
    "parameters": ......
}
```

Figure 20: An example of Special Irrelevant Data.

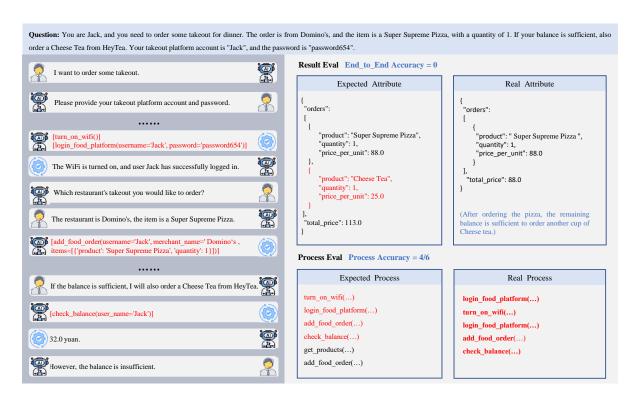


Figure 21: An example of Agent Data: A partially correct representation.

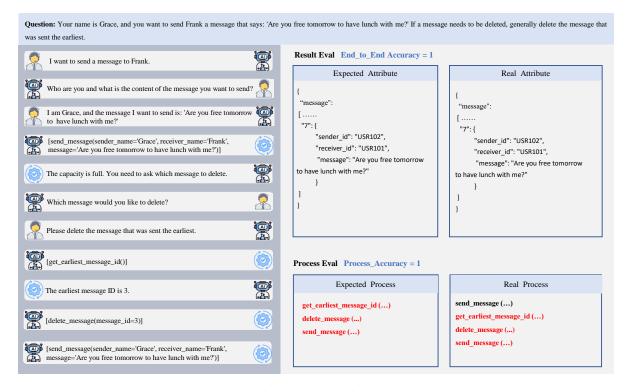


Figure 22: An example of Agent Data: A fully accurate representation.

You are an AI assistant with the role name "assistant". Based on the provided API specifications and conversation history from steps 1 to t, generate the API requests that the assistant should call in step t+1. The API requests should be output in the format [ApiName(key1='value1', key2='value2', ...)], replacing ApiName with the actual API name, key1, key2, etc., with the actual parameter names, and value1, value2, etc., with the actual parameter values. The output should start with a square bracket "[" and end with a square bracket "]". If there are multiple API requests, separate them with commas, for example:[ApiName(key1='value1',key2='value2',...), ApiName(key1='value1',key2='value2',...), ...]. Do not include any other explanations, prompts, or API call results in the output. If the API parameter description does not specify otherwise, the parameter is optional (parameters mentioned in the user input need to be included in the output; if not mentioned, they do not need to be included). If the API parameter description does not specify the required format for the value, use the user's original text for the parameter value. If the API requires no parameters, output the API request directly in the format [ApiName()], and do not invent any nonexistent parameter names.

{time}

**Role Descriptions:** 

user: User

assistant: The AI assistant role that makes API requests tool: Provides the results returned from tool calls

**API Specifications:** 

{function}

Figure 23: The inference prompt for Nomal (except Prefernce) Data.

You are an AI assistant, and your role is called assistant. **Based on the given API description, dialogue history 1..t, and character profile**, generate the API requests that the assistant should call in step t+1. The API requests should be output in the format [ApiName(key1='value1', key2='value2', ...)], where ApiName is replaced with the actual API name, and key1, key2, etc., are replaced with the actual parameter names, and value1, value2 are replaced with the actual parameter values. The output should start with a "[" and end with a "]". If there are multiple API requests, they should be separated by commas, e.g., [ApiName(key1='value1', key2='value2', ...), ApiName(key1='value1', key2='value2', ...), ...]. Do not output any other explanations, hints, or results of the API calls in the output. If the API parameter description does not specify special instructions, the parameter is optional (parameters mentioned in the user input or character profile should be included in the output, and if not mentioned, they should not be included). If the API parameter description does not specify the format for the parameter value, it should be taken from the user's original text or character profile. If the API requires no parameters, the API request should be output as [ApiName()], with no fabricated parameter names.

Character Profile:

{profile}

**Role Descriptions:** 

user: User

assistant: The AI assistant role that makes API requests tool: Provides the results returned from tool calls

**API Specifications:** 

{function}

Figure 24: The inference prompt for Nomal (Prefernce) Data.

You are an AI assistant with the role name "assistant". Based on the provided API specifications and conversation history from steps 1 to t, generate the API requests that the assistant should call in step t+1. The API requests should be output in the format [ApiName(key1='value1', key2='value2', ...)], replacing ApiName with the actual API name, key1, key2, etc., with the actual parameter names, and value1, value2, etc., with the actual parameter values. The output should start with a square bracket "[" and end with a square bracket "]".

#### {time}

#### **Role Descriptions:**

user: User

assistant: The AI assistant role that makes API requests tool: Provides the results returned from tool calls

#### **API Specifications:**

{function}

Figure 25: Condensed Prompt for Normal Data.

You are an AI assistant. Based on the provided API specifications and conversation history generate the API requests in the format [ApiName(key1='value1', key2='value2', ...),.....].

#### {time}

#### **Role Descriptions:**

user: User

assistant: The AI assistant role that makes API requests tool: Provides the results returned from tool calls

#### **API Specifications:**

{function}

Figure 26: Minimal Prompt for Normal Data.

You are an AI system with the role name "assistant". Based on the provided API specifications and conversation history from steps 1 to t, generate the API requests that the system should call in step t+1. Below are two specific scenarios:

- 1. When the information provided by the user is clear and unambiguous, and the problem can be resolved using the list of candidate functions:
- If the API parameter description does not specify the required format for the value, use the user's original text for the parameter value.
- API requests should be output in the format [ApiName(key1='value1', key2='value2', ...), ApiName(key1='value1', key2='value2', ...), ...], replacing ApiName with the actual API name, key1, key2, etc., with the actual parameter names, and value1, value2, etc., with the actual parameter values. The output should start with a square bracket "[" and end with a square bracket "]". At this time, the output must not contain any other content.
- 2. When the information provided by the user is unclear, incomplete, or incorrect, or the user's question exceeds the capabilities of the provided functions, you need to clearly point out these issues. The following is your strategy: (1) If the user's instructions include the key details required to call the API, but the type or form of the parameter values does not match the API's definitions, ask in-depth questions to clarify and correct the details. The output format should be: ["There is incorrect value (value) for the parameters (key) in the conversation history."]
- (2) If the user's instructions lack the key details required by the API, ask questions to obtain the necessary information. The output format should be: ["Missing necessary parameters (key1, key2, ...) for the api (ApiName)"], replacing key1, key2 with the names of the missing parameters and ApiName with the actual API name.
- (3) If the user's request exceeds the current capabilities of your APIs, inform them that you cannot fulfill the request. The output format should be: ["Due to the limitations of the function, I cannot solve this problem."]

Note: The above steps have a priority order. You need to first determine whether scenario (1) applies. If it does, output according to the requirements in (1). Pay attention to distinguishing between scenarios (1) and (2).

{time}

#### **Role Descriptions:**

user: User

assistant: The AI assistant role that makes API requests tool: Provides the results returned from tool calls

#### **API Specifications:**

{function}

Figure 27: The inference prompt for Special Data.

The current time is June 11, 2024, 16:00 (Beijing Time). As a simulated mobile assistant agent, you can help users **send text messages, add reminders, and order takeout.** 

#### Text messages

**Sending Text Messages** (1)Before sending a text message, the agent must first obtain the sender and recipient of the message.(2)When the memory is full and needs to delete messages, you need to ask the user: "Memory is full, which message would you like to delete?"

**Viewing Text Messages** (1)Before viewing text messages, the agent must first log into the device via login\_device().(2)Before viewing text messages, the agent must first obtain the sender and recipient of the messages.(3)After viewing text messages, the agent needs to ask the user if they want to add the message content to a reminder.(4)After viewing text messages, the agent needs to ask the user if they want to reply to the message.(5)If the message content involves takeout, the agent needs to ask if the user wants to order takeout based on the message content.

#### Reminders

**Adding Reminders**(1)Before adding a reminder, you should obtain the content and title of the reminder. The reminder time defaults to the current time.(2)If the reminder to be added is the content of a specific message, the agent needs to first view the message content.

Viewing Specific Reminders by Title: After viewing a specific reminder by title, you need to ask the user if they want to complete the tasks within it.

#### Order takeout

Ordering Takeout(1)Before ordering takeout, the agent needs to obtain the user's takeout platform account and password, and log in using login\_food\_platform().(2)If the merchant, product, and quantity for the order are not initially provided, you need to ask the user.(3)When encountering takeout from different merchants, you need to order them one by one.(4)If the balance is insufficient, you need to inform the user "Insufficient balance" and ask if they want to change the order.

#### **Function Calls**

When a function call is needed, please strictly adhere to the above format requirements:

(1)[ApiName(key1='value1', key2='value2', ...)], Please remember that the function call must start with [ and end with ] (2)You need to promptly feedback the task execution status to the user and do not repeatedly call the same function. When you believe the current task is completed, respond with "finish conversation" to end the dialogue.

Figure 28: The inference prompt for Agent Data in a specific scenario.

As a user, your role is to interact with an agent. However, during the interaction, you need to follow these guidelines:

- 1 Break down your inquiries and only raise one question per exchange to simulate a real user's messages.
- 2 Provide all the necessary information for the current step. For instance, when setting a reminder, you must give details such as the reminder's description, title, and time.
- 3 When asked if you require further assistance, ensure that the main tasks in the instruction have been completed. If not, continue to present the next step to the agent.
- 4 When the agent asks which message needs to be deleted, proceed with the deletion as specified in the instructions. You cannot offer proactive help to the agent; respond to the agent's questions according to the instructions, and do not invent any information that you do not know.
- 5 Once all tasks are complete, generate a 'finish conversation' message as a standalone line to end the discussion. Question: question

Figure 29: An inference prompt for user simulator in a specific scenario.

Table 7: Accuracy evaluation of different models on English Data (%).

Model			N	ormal			Special	Agent	Overall
	Atom	Single-Turn	Multi-Turn	Similar API	Preference	Summary	Special		
		Closed	-Source Large	Language Mo	dels				
GPT-4o	90.0	78.0	68.0	80.0	78.0	82.5	92.7	56.0	81.1
GPT-4-Turbo	90.7	80.5	69.0	80.0	88.0	84.2	82.0	62.5	80.3
Qwen-Max	88.0	75.0	61.0	74.0	82.0	79.7	74.0	60.0	75.1
GPT-4o-Mini	84.3	73.5	59.0	74.0	72.0	76.4	76.7	27.5	68.9
Gemini-1.5-Pro	82.3	73.0	61.0	74.0	72.0	75.7	77.3	26.0	68.5
Claude-3-5-Sonnet	66.7	64.0	46.0	58.0	68.0	62.2	72.7	44.0	62.2
Doubao-Pro-32k	75.3	58.0	52.0	70.0	54.0	66.3	50.7	26.5	56.0
		Open-	Source Large	Language Mod	dels				
Qwen2.5-Coder-32B-Instruct	86.0	73.5	59.0	76.0	72.0	77.4	80.0	50.0	73.9
DeepSeek-V3	88.0	77.5	63.0	76.0	78.0	80.3	72.7	34.0	71.1
Qwen2.5-72B-Instruct	81.3	74.5	64.0	76.0	80.0	76.8	74.0	37.5	70.0
Llama-3.1-70B-Instruct	83.7	71.5	61.0	74.0	66.0	75.6	29.3	41.0	57.9
Qwen2.5-7B-Instruct	70.3	57.0	49.0	62.0	58.0	62.8	49.3	15.0	51.8
Qwen2.5-Coder-7B-Instruct	73.3	63.5	52.0	70.0	58.0	66.6	25.3	18.5	48.1
DeepSeek-Coder-V2-Lite-Instruct	71.7	58.0	50.0	62.0	60.0	64.0	39.3	2.5	47.9
Watt-Tool-8B	84.7	71.5	57.0	70.0	62.0	74.8	2.0	1.5	44.0
Hammer2.1-7B	71.3	62.5	43.0	64.0	52.0	62.9	3.3	15.0	39.6
Phi-3-Mini-128k-Instruct	66.3	49.0	31.0	58.0	32.0	54.0	12.0	0.0	34.4
MLlama-3.1-8B-Instruct	51.0	49.5	28.0	60.0	56.0	48.1	15.3	6.5	32.9
xLAM-7B-r	61.7	42.0	32.0	66.0	0.0	48.7	4.0	10.0	30.8
Llama-3.2-3B-Instruct	31.7	21.5	9.0	34.0	32.0	26.4	8.7	0.0	17.6
Hammer2.1-3B	32.7	14.0	7.0	36.0	32.0	25.5	0.7	1.5	15.2

Table 8: Accuracy evaluation of different models on Chinese Data (%).

Model	Normal							Agent	Overall
Wiodei	Atom	Single-Turn	Multi-Turn	Similar API	Preference	Summary	Special	rigent	Overun
		Closed	-Source Large	Language Mo	dels				
GPT-4o	96.7	91.0	86.0	90.0	88.0	92.7	93.3	71.5	89.6
GPT-4-Turbo	95.7	89.0	86.0	92.0	84.0	91.7	91.3	72.5	88.6
Qwen-Max	94.3	86.0	75.0	92.0	84.0	88.7	74.0	68.5	81.7
GPT-4o-Mini	88.7	78.5	74.0	80.0	84.0	83.4	81.3	39.0	76.0
Claude-3-5-Sonnet	87.0	81.0	79.0	84.0	76.0	83.5	82.0	35.0	75.6
Gemini-1.5-Pro	86.7	80.5	68.0	86.0	84.0	82.2	80.0	25.0	72.8
Doubao-Pro-32k	84.3	53.0	64.0	82.0	78.0	75.0	59.3	23.5	62.8
		Open-	Source Large	Language Mod	dels				
Qwen2.5-Coder-32B-Instruct	94.3	88.5	83.0	90.0	90.0	90.8	81.3	71.5	85.3
Qwen2.5-72B-Instruct	92.3	86.0	75.0	90.0	82.0	87.3	77.3	52.5	79.3
DeepSeek-V3	95.0	90.5	91.0	90.0	88.0	92.6	73.3	35.0	78.5
Llama-3.1-70B-Instruct	81.3	65.0	66.0	84.0	70.0	75.3	47.3	43.5	62.9
Qwen2.5-7B-Instruct	81.7	63.5	68.0	82.0	76.0	75.9	44.7	12.5	57.8
DeepSeek-Coder-V2-Lite-Instruct	78.7	57.5	43.0	82.0	70.0	68.8	41.3	1.5	51.1
Qwen2.5-Coder-7B-Instruct	78.7	64.0	63.0	78.0	78.0	73.5	19.3	12.5	49.6
Watt-Tool-8B	86.7	67.0	54.0	88.0	66.0	76.3	10.0	4.0	47.4
Hammer2.1-7B	76.0	62.5	37.0	60.0	58.0	62.7	26.0	18.5	46.1
Llama-3.1-8B-Instruct	52.7	30.0	28.0	72.0	36.0	45.0	26.7	4.0	33.8
Phi-3-Mini-128k-Instruct	48.0	29.5	15.0	58.0	32.0	38.9	25.3	1.5	29.5
Llama-3.2-3B-Instruct	45.7	9.0	9.0	50.0	32.0	32.7	10.0	0.0	21.6
xLAM-7B-r	25.3	2.0	6.0	56.0	0.0	18.7	1.3	7.5	12.3
Hammer2.1-3B	12.0	9.0	0.0	44.0	8.0	11.8	1.3	1.5	7.4

```
(question)
I want to understand the symmetry in Escher's Waterfallärtwork. Please provide a detailed symmetry analysis.

(answer)
{"Escher_Artwork_Analysis_Tool": {
    "artwork_analysis": [
        {"artwork_id": "Waterfall",
        "analysis_type": "Symmetry Analysis"}]}}

(model output)
{"FontAdvisor_recommendDisplayFonts":{......}}
```

Figure 30: An error example of Wrong Function Name on Normal Data.

```
(question)
Can you check if our current automation system with system ID 2435 is compatible with the new device model XR10?

(answer)
{"AutomationSystemCompatibilityCheck":{ "systemID":2435,"newDeviceModel":"XR10"}}

(model output)
{"AutomationSystemCompatibilityCheck":{ "systemID":"2435", "newDeviceModel":"XR10"}}
```

Figure 31: An error example of Wrong Param Type on Normal Data.

```
(question)
I am currently working on a research paper about the state of forests, and I need an assessment of the vegetation diversity
in the Amazon rainforest. Please provide me with a detailed report.

(answer)
{"EcoDiversityMonitor_assessVegetation": {
    "area_identifier": "Amazon rainforest",
    "detailed_report": true}}
(model output)

[user,area_identifier='Amazonrainforest',detailed_report=True]
```

Figure 33: An error example of Wrong Output Format on Normal Data.

```
(question)
I would like to assess the GPS accuracy of my vehicle. The device ID is GPS-7524. Sorry, I made a mistake earlier—it should be GPS-7526.

(answer)
{ "SignalAccuracyEvaluator":{ "device_id":"GPS-7526" }}

(model output)
{ "SignalAccuracyEvaluator":{ "device_id":"GPS-7524"}}
```

Figure 32: An error example of Wrong Param Value on Normal Data.

```
(question)
I'm considering relocating my business to the Middle East. Can you provide me with a detailed list of major cities there, sorted by internet speed?

(answer)
Missing necessary parameters (sort) for the API (Get_Middle_East_Cities)

(model output)
{
    "Get_Middle_East_Cities": {
        "sort": "desc",
        "sort_by": "internet_speed",
        "size": 10,
        "page": 1
    }
}
(The question did not mention the sort method)
```

Figure 34: An error example of Error Detection on Special Data.

## (question) I'm considering relocating my business to the Middle East. Can you provide me with a detailed list of major cities there, sorted by internet speed? (answer) Missing necessary parameters (sort) for the api (Get\_Middle\_East\_Cities) (model output) Missing necessary parameters (sort by) for the api (Get\_Middle\_East\_Cities) (Find missing necessary parameters, but not correctly identified)

Figure 35: An error example of Error Correction on Special Data.