

ASKToACT: Enhancing LLMs Tool Use via Self-Correcting Clarification

Anonymous ACL submission

Abstract

Large language models (LLMs) have demonstrated remarkable capabilities in tool learning. In real-world scenarios, user queries are often ambiguous and incomplete, requiring effective clarification. However, existing interactive clarification approaches face two critical limitations: reliance on manually constructed datasets, which inherently constrains training data scale and diversity, and lack of error correction mechanisms during multi-turn clarification, leading to error accumulation that compromises both accuracy and efficiency. We present **ASKToACT**, which addresses these challenges by exploiting the structural mapping between queries and their tool invocation solutions. Our key insight is that tool parameters naturally represent explicit user intents. By systematically removing key parameters from queries while retaining them as ground truth, we enable automated construction of high-quality training data. We further enhance model robustness through error-correction pairs and selective masking, enabling dynamic error detection during clarification interactions. Comprehensive experiments demonstrate that **ASKToACT** significantly outperforms existing approaches, achieving above 57% accuracy in recovering critical unspecified intents and enhancing clarification efficiency by an average of 10.46% while maintaining high accuracy in tool invocation. Our framework exhibits robust performance across different model architectures and successfully generalizes to entirely unseen APIs without additional training, achieving performance comparable to GPT-4o with substantially fewer computational resources.

1 Introduction

Large language models (LLMs) have demonstrated remarkable capabilities in various tasks, from code generation to complex reasoning (Nakano et al., 2021; Chen et al., 2021; Komeili et al., 2022; Wei

et al., 2022). A particularly promising direction is their ability to interact with external tools through API calls, which significantly expands their practical applications (Schick et al., 2023; Hao et al., 2023; Qin et al., 2024; Shim et al., 2025). This has inspired numerous frameworks focusing on tool-augmented LLMs, including Toolformer (Schick et al., 2023), ToolLLaMA (Qin et al., 2024), and Gorilla (Patil et al., 2023).

However, current tool learning frameworks (Li et al., 2023; Song et al., 2023; Schick et al., 2023; Qin et al., 2024) operate under an idealistic assumption that user queries are always explicit and unambiguous. This diverges significantly from real-world scenarios where users often provide incomplete, ambiguous, or imprecise queries. Such unspecified queries pose unique challenges in tool learning scenarios, as API calls require precise parameters and cannot tolerate ambiguity (Wang et al., 2024b). When faced with unspecified queries, LLMs tend to either arbitrarily generate missing parameters or remain unknown, leading to potential risks in tool invocation.

This raises a critical research question: *How can we enhance LLMs’ ability to handle unspecified queries in tool learning scenarios while ensuring accurate and reliable tool invocation?* Recent works (Zhang and Choi, 2023; Qian et al., 2024; Wang et al., 2024b) have introduced interactive clarification approaches, but face two fundamental limitations. First, they rely heavily on manually constructed datasets for training (Qian et al., 2024; Wang et al., 2024b). Creating these datasets requires human annotators to craft queries and clarifications, a process that inherently limits scale and diversity. The resulting datasets capture only a narrow range of ambiguity patterns, reducing their effectiveness with diverse real-world queries. Second, as shown in Figure 1, these approaches lack robust error handling during multi-turn clarification. Existing models train on datasets with only

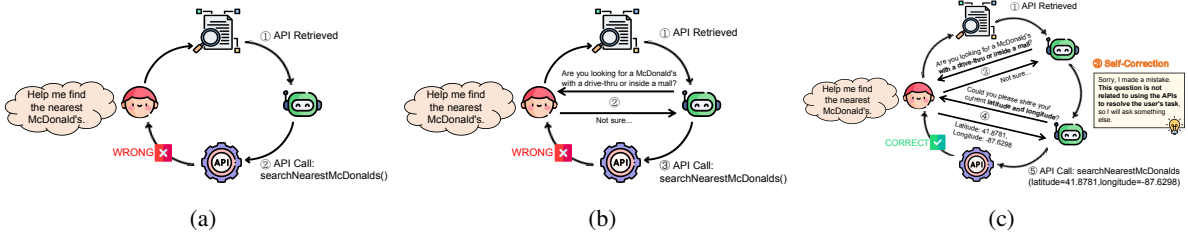


Figure 1: Comparison of query handling approaches: (a) direct API calls without clarification, (b) basic clarification without error recovery, and (c) our self-correcting **ASKTOACT** framework.

perfect clarification sequences. In reality, models often request already-provided information, follow irrelevant paths, or miss unspecified details. Without error recovery training, these issues accumulate throughout dialogues, reducing efficiency and compromising tool invocation quality.

We propose **ASKTOACT**, a self-correcting clarification framework that systematically addresses these limitations. Our key insight is that tool parameters naturally represent explicit user intents, creating an opportunity for automated data generation. We develop an automated pipeline that strategically removes key parameters from complete queries in existing datasets, generating diverse unspecified queries with built-in ground truth. Using these queries, we construct rich clarification dialogues demonstrating effective intent elicitation. To enable robust error handling during interactions, we augment training data with carefully designed error-correction pairs that simulate realistic mistakes and their solutions. We implement selective masking during training to prevent learning negative patterns while enhancing error detection abilities.

Through comprehensive experiments, we demonstrate that **ASKTOACT** achieves several significant improvements: (1) correctly identifies unspecified queries and recovers more than 57% of critical unspecified intents, while significantly enhances clarification efficiency by an average of 10.46% compared to the base model; (2) achieves strong performance in end-to-end tool invocation, with over 81% tool selection accuracy and over 68% parameter resolution accuracy; (3) exhibits robust performance across different model architectures, and successfully generalizes to entirely unseen APIs; and (4) delivers performance comparable to GPT-4o while requiring substantially fewer computational resources.

Our work makes three main contributions:

- We introduce an automated pipeline for generating high-quality intent clarification datasets, addressing the scalability limitations of manual annotation.
- We develop a self-correction mechanism that enables models to dynamically detect and correct potential errors during clarification interactions.
- Our experimental results demonstrate that our method not only achieves state-of-the-art (SOTA) performance but also shows strong generalization ability when handling queries requiring the use of unseen APIs.

2 Method

Tool learning faces a fundamental challenge: while API calls require precise parameters, real-world queries are often ambiguous. To bridge this gap, we propose **ASKTOACT**, a self-correcting clarification framework. Our method consists of two key components: (1) an automated data construction pipeline for generating diverse intent clarification data (§2.1), and (2) a self-correction training paradigm for dynamic error detection and correction (§2.2). The core insight is that tool parameters naturally represent explicit user intents, making them ideal anchors for both data generation and error correction. Figure 2 illustrates the overall framework architecture.

2.1 Intent Clarification Dataset Curation

The foundation of our method is a systematic pipeline for constructing multi-turn clarification data. As shown in Figure 2, the pipeline proceeds in two steps: generating unspecified queries and subsequently constructing clarification dialogues.

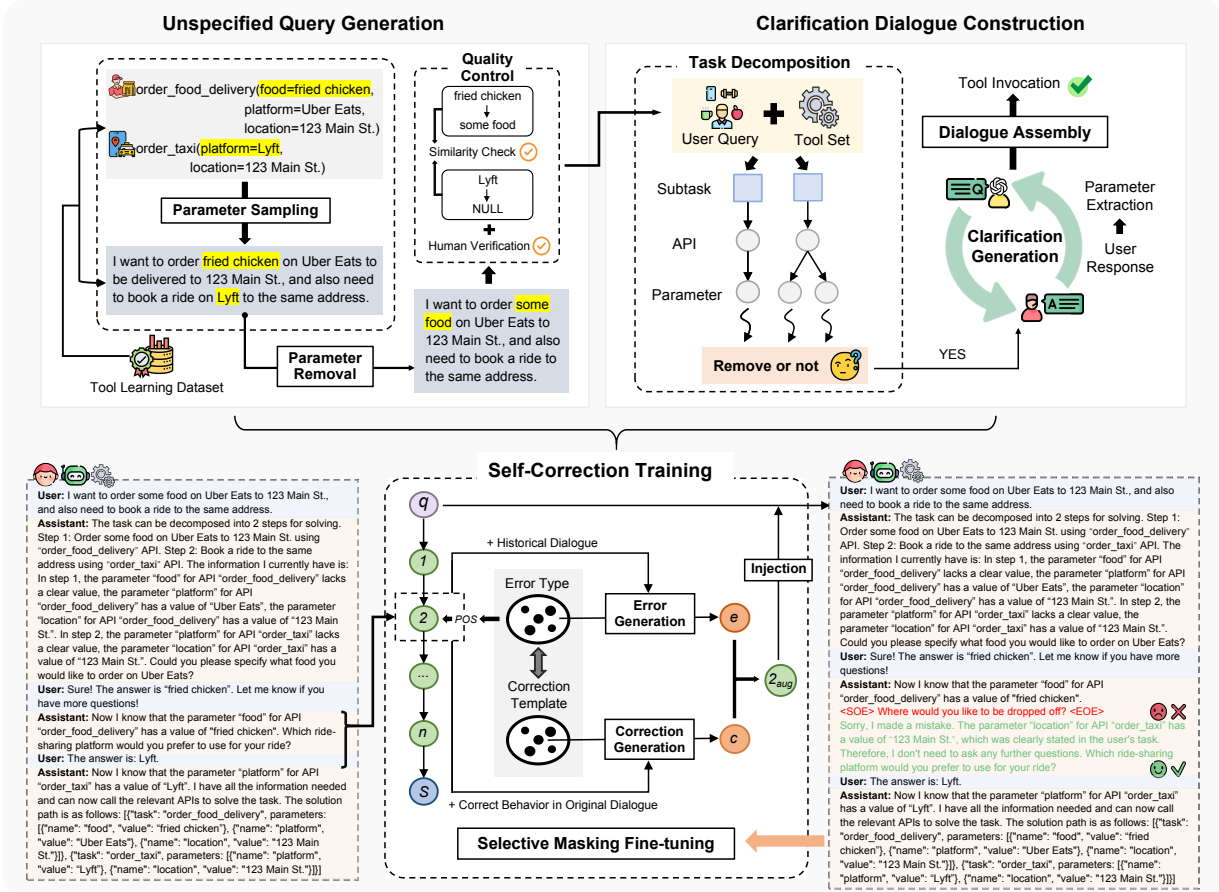


Figure 2: Overview of ASKTOACT framework. Top: Dataset construction pipeline, consisting of (1) unspecified query generation by selecting and removing key parameters (e.g., "fried chicken" and "Lyft") from original queries, and (2) dialogue construction through task decomposition, clarification generation, and dialogue assembly. Bottom: Self-correction training through dialogue augmentation and selective masking fine-tuning.

2.1.1 Unspecified Query Generation

A key challenge in building intent clarification systems is obtaining realistic examples of ambiguous queries paired with their complete intents. We address this through a novel reverse-engineering approach that leverages existing tool learning datasets. Each instance in these datasets consists of a fully specified query q and a corresponding tool invocation solution $S = \{(f_i, P_i) \mid i = 1, \dots, n\}$, where f_i denotes the API and $P_i = \{p_i^1, p_i^2, \dots\}$ represents its parameter set. We systematically transform q into an unspecified query q' while preserving the ground truth information necessary for subsequent dialogue construction and evaluation.

Parameter Sampling The first step in our pipeline is to determine which parameters in S to remove from the original query q . We implement a stratified sampling approach that enables sampling across different API domains and parameter counts. For each query, we randomly select

parameters according to one of four complexity levels: (1) fully specified, where all parameters are retained, (2) single-API single-parameter, where one parameter from one API call is removed, (3) single-API multi-parameter, where multiple parameters from the same API call are removed, and (4) multi-API multi-parameter, where parameters are removed across multiple API calls. This stratification ensures our dataset captures the full spectrum of query ambiguity encountered in real-world scenarios, from basic single-slot ambiguities to complex multi-faceted ambiguities.

Parameter Removal Once the removed parameter set P' is determined, we apply two complementary strategies to transform the original query q into an unspecified form. The first strategy, *complete removal*, entirely eliminates parameter values from q while preserving grammatical integrity. The second strategy, *semantic abstraction*, replaces specific parameter values with abstract expressions that ne-

cessitate further clarification. For each parameter $p \in P'$, we maintain a mapping $M : p \rightarrow v'$, where v' represents the transformed value after parameter removal or abstraction. By recording the values before and after the transformation, we can more precisely track how explicit user intent becomes unspecified during the process of unspecified query generation. This record plays an important role in quality control, helping us ensure the quality of the generated query. Also, it guides subsequent dialogue construction and provides ground truth for evaluation. The implementation details and the format of the transformation record are provided in Appendix C.1.1 and Appendix C.1.2, respectively.

Quality Control To ensure generation quality, we employ a dual-stage verification mechanism. We first compute semantic similarity between original values and their transformations using Sentence Transformer (Reimers and Gurevych, 2019). Queries with similarity scores exceeding 0.95 are filtered out to ensure sufficient semantic alteration. Then, we conduct human verification on generated queries (see Appendix C.1.3). Through this process, we constructed 35,261 high-quality unspecified queries, as shown in Table 1.

2.1.2 Clarification Dialogue Construction

Based on the generated unspecified queries, we propose an automated method to construct training dialogue data that simulate multi-turn clarification. The dialogue construction process—from task decomposition, clarification generation to final dialogue assembly—is essential for generating coherent and effective clarification interactions.

Task Decomposition The foundation of effective clarification lies in identifying what information needs to be clarified. Given an unspecified query q' and its tool invocation solution S , we first decompose the query into a sequence of subtasks. Each subtask corresponds to an API call in S . For each API call, we perform parameter analysis to identify two categories: (1) parameters already specified in q' , and (2) parameters requiring clarification. This structured decomposition guides the subsequent clarification process, ensuring all necessary information is systematically obtained.

Clarification Generation Building on the task decomposition, we generate interaction turns for each parameter that requires clarification, following the API call order defined in S . We construct

each clarification turn through a three-step process, with the goal of maximizing clarification effectiveness and preserving natural conversational flow. First, we generate a clarification question q_c targeting the unspecified parameter. Next, we simulate the user’s reply using diverse response templates that vary in verbosity and conversational tone. Finally, we generate a confirmation statement summarizing the clarified information, which serves as an explicit reference for subsequent turns.

Dialogue Assembly The final step brings together all components into a coherent dialogue structure. We assemble the generated elements sequentially while maintaining natural conversation flow through consistent reference to previously clarified information and smooth transitions between parameter-related clarifications. Special attention is paid to parameter interdependencies, ensuring that information is requested in a logical order that reflects real-world dialogue patterns. The dialogue concludes with the complete tool invocation solution S , providing a clear connection between the clarification process and its ultimate goal. Detailed templates and prompting strategies that support this assembly process are provided in Appendix C.2.

2.2 Self-Correction Training

While constructing high-quality training data is essential, the dynamic nature of clarification interactions requires models to detect and correct potential errors in real-time. We develop a systematic training paradigm that combines error-correction augmentation with specialized training strategies to enhance model robustness and enable self-correction.

Error Type Analysis Through comprehensive analysis of clarification interactions generated by the model in response to unspecified queries, we identify five primary error types that impair the clarification process in complementary ways. *Clearly Stated Intent Clarification* occurs when the model requests explicitly stated information, creating unnecessary interaction turns. *Imprecise Clarification* is characterized by questions that lack specificity, often resulting in ambiguous user responses. *Irrelevant Clarification* emerges when the model poses questions that diverge from the core intent. *Redundant Clarification* arises when the model requests information that has been previously clarified. *Incomplete Clarification* represents failure to identify all parameters requiring clarification, leading to incomplete tool invocation solutions. Understanding

| Dataset | Train | Test | Avg. No. APIs | Avg. No. Params | Avg. No. Unspecified Intents |
|--------------|--------|-------|---------------|-----------------|------------------------------|
| xlam-IC | 29,821 | 4,456 | 1.58 | 2.49 | 1.32 |
| Taskbench-IC | - | 984 | 1.75 | 2.35 | 1.30 |
| Total | 29,821 | 5,440 | 1.59 | 2.49 | 1.32 |

Table 1: Datasets statistics. xlam-IC is generated from the xlam-function-calling-60k dataset (Liu et al., 2024) which is used for training and testing. Taskbench-IC is generated from Taskbench (Shen et al., 2024) and is used exclusively for OOD testing. Please refer to the Appendix B for more details about both datasets.

| Error Type | Count |
|-------------------------------------|--------|
| Clearly Stated Intent Clarification | 2,481 |
| Imprecise Clarification | 2,298 |
| Irrelevant Clarification | 2,251 |
| Redundant Clarification | 3,126 |
| Incomplete Clarification | 5,000 |
| Total | 15,156 |

Table 2: Error type distribution in augmented dialogues.

these patterns guides our error-correction strategy.

2.2.1 Error-Correction Augmentation

Building upon this error analysis, we introduce an automated method to augment dialogues with error-correction pairs. Given a dialogue d , we randomly select an error type τ_k and determine an injection position pos . We then generate the error instance using two strategies. For semantic errors (Clearly Stated Intent, Imprecise, and Irrelevant), we employ GPT-4o with specialized prompts: $e = f_{\text{gpt}}(d, \tau_k, pos)$. For structural errors (Redundant and Incomplete), we implement rule-based algorithms: $e = f_{\text{rule}}(d, \tau_k, pos)$. Implementation details are provided in Appendix D.1.

After generating error e , we construct a correction c using a template specific to the error type τ_k . The resulting correction c explicitly states the error type and identifies the correct behavior as the content at position pos in the original dialogue d . We then inject these error-correction pair (e, c) into the original dialogue d at the predetermined position pos : $d' = \text{inject}(d, e, c, pos)$.

To ensure the validity of our error-correction augmentation method, we conduct human verification on augmented dialogues (see Appendix D.3). Through this systematic process, we generated 15,156 augmented dialogues, as shown in Table 2.

2.2.2 Selective Masking Fine-tuning

To effectively utilize the augmented dialogues for training while preventing the model from learning error patterns, we implement a selective mask-

ing mechanism during fine-tuning. We introduce special tokens $\langle SOE \rangle$ and $\langle EOE \rangle$ to demarcate error segments, and mask these segments during loss computation. This approach allows the model to learn error detection and correction patterns while avoiding the reinforcement of error behaviors. Through this training process, we equip the model with the ability to dynamically identify potential errors and apply appropriate corrections during clarification interactions.

3 Experiment

3.1 Experimental Settings

Training Details We construct our training data from the xlam-IC dialogue dataset, where 30% of the samples are replaced with error-correction augmented dialogues. We explore two adaptation strategies for the Qwen2.5-7B-Instruct model: LoRA (Hu et al., 2021) and full-parameter fine-tuning. More details are provided in Appendix E.

Baselines For comprehensive comparison, we evaluate representative tool-augmented LLMs, including xLAM-7b-fc-r (trained on the xlam-function-calling-60k dataset but without intent clarification) (Liu et al., 2024), gorilla-openfunctions-v2 (Patil et al., 2023), and ToolLLaMA-2-7b-v2 (Qin et al., 2024), as well as an intent clarification model, Mistral-Interact (Qian et al., 2024). In addition, we evaluate major LLM series, including Mistral-7B-Instruct-v0.3, LLaMA (3-8B/70B-Instruct), Qwen (2.5-7B/72B-Instruct), DeepSeek-V3, Claude (3.5-Haiku/Sonnet), and GPT (3.5, 4, 4o-mini, 4o). All models use a standardized evaluation prompt (see Appendix F.1).

3.2 Evaluation Framework

We develop an automated framework for systematic evaluation on handling unspecified queries. The framework employs an LLM to simulate user behavior. During interactions, the user-simulating LLM judges whether clarification questions are relevant to unspecified intents, and either provides the

necessary information or indicates that it is unavailable. To better capture the complexity of real-world human-LLM interactions, we configure the user-simulating LLM with six personality types, each exhibiting different response characteristics. Implementation details are provided in Appendix F.2.

3.3 Metrics

We evaluate the models in two aspects: *intent clarification quality* and *tool invocation accuracy*. For intent clarification quality, we design four metrics. *Intent Coverage Rate (ICR)* measures the proportion of successfully clarified intents among all unspecified intents, while *Clarification Efficiency (CE)* evaluates the success rate of clarification across interaction rounds. We combine these measures into a *Clarification Performance Score (CPS)* using a harmonic mean, similar to the F1-score formulation. Additionally, we track *Interaction Rounds (IR)* as the average number of clarification rounds per query. For tool invocation accuracy, we introduce three complementary metrics. *Solution Completion Rate (SCR)* measures the proportion of successfully generated tool invocation solutions, providing an end-to-end assessment. *Tool Selection Score (TSS)* evaluates API selection accuracy using an F1-score over selected and required APIs. *Parameter Resolution Score (PRS)* assesses the accuracy of parameter resolution through an F1-score computation over API-parameter-value triples. The details are provided in Appendix F.3.

3.4 Main Results

3.4.1 LLM-based Simulated Evaluation

The experimental results on the in-domain (ID) test split of the xlam-IC dataset are presented in Table 3. Our method demonstrates superior performance in both intent clarification and tool invocation.

Intent Clarification Capability Both variants of our method—ASKToACT-LoRA-SFT-7B and ASKToACT-Full-SFT-7B—exhibits strong capabilities in intent clarification. In particular, the fully fine-tuned variant reaches a CPS of 65.92%, closely approaching the performance of SOTA LLMs such as GPT-4o. Meanwhile, the lightweight LoRA variant also achieves competitive results (ICR: 57.68%, CE: 63.41%, CPS: 60.41%), significantly surpassing the specialized intent clarification model Mistral-Interact.

Tool Invocation Accuracy Our method demonstrates remarkable capabilities in translating clar-

ified intents into precise tool invocations. It achieves SOTA performance across all evaluation metrics (SCR > 96%, TSS > 81%, PRS > 68%), significantly surpasses all existing open-source and closed-source LLMs. Compared to tool-augmented LLMs such as xLAM-7b-fc-r, gorilla-openfunctions-v2, and ToolLLaMA-2-7b-v2, our method demonstrates substantial advantages. This performance gap highlights the strength of integrating intent clarification with tool learning. Unlike prior specialized models that are limited to unambiguous tool-use queries, our method effectively resolves ambiguity in user queries, leading to significantly improved tool invocation accuracy.

Further analyses—including cross-model transferability, the impact of augmentation proportion, clarification complexity, and a case study of user interaction styles—are presented in Appendix G.

3.4.2 Human-Interactive Evaluation

To assess the effectiveness of our method in real-world interactions, we conducted a human-interactive evaluation. We recruited 3 participants, each asked to propose 10 unspecified tool-use queries requiring clarification. These queries were independently tested on the base model (Qwen2.5-7B-Instruct) and our models (ASKToACT-LoRA-SFT-7B and ASKToACT-Full-SFT-7B). Participants interacted with the models iteratively until they obtained a satisfactory response.

As shown in Table 4, both variants of our method outperform the base model across all metrics. Specifically, ASKToACT-LoRA-SFT-7B improves the Task Completion Rate by 6.66% and the Intent Coverage Rate by 9.76%, while reducing Interaction Rounds from 3.20 to 2.73. ASKToACT-Full-SFT-7B achieves further improvements, reaching 96.67% Task Completion Rate and 85.37% Intent Coverage Rate, with fewer Interaction Rounds (2.60). In addition, participants reported higher satisfaction with both variants (4.40 and 4.61 vs. 3.80), confirming that our method leads to a more effective and user-friendly interaction experience. The consistency of results across both LLM-based and human-interactive evaluations highlights the effectiveness, robustness, and practical utility of our method.

3.5 OOD Generalization

To assess the generalization ability of our method, we test on Taskbench-IC, an out-of-domain (OOD) set that consists of entirely unseen API domains.

| LLM | Intent Clarification Quality | | | | Tool Invocation Accuracy | | |
|---------------------------|------------------------------|----------------|-----------------------|--------------|--------------------------|-----------------------|-----------------------|
| | ICR↑ | CE↑ | CPS↑ | IR↓ | SCR↑ | TSS↑ | PRS↑ |
| <i>Closed-Source LLMs</i> | | | | | | | |
| Claude3.5-Haiku | 49.60 | 35.05 | 41.07 | 2.30 | 84.20 | 62.74 | 52.12 |
| Claude3.5-Sonnet | 57.55 | 61.71 | 59.55 | 1.52 | 94.52 | 73.20 | 62.68 |
| GPT-3.5 | 46.63 | 51.41 | 48.90 | 1.48 | 93.20 | 67.75 | 51.22 |
| GPT-4 | 59.43 | 63.09 | 61.21 | 1.53 | 93.42 | 71.55 | 61.82 |
| GPT-4o-Mini | 57.95 | 56.43 | 57.18 | 1.67 | 92.98 | 71.82 | 61.52 |
| GPT-4o | 64.82 | 74.50 | 69.33 | 1.33 | 94.52 | 76.94 | 67.65 |
| <i>Open-Source LLMs</i> | | | | | | | |
| Mistral-7B-Instruct-v0.3 | 26.01 | 34.90 | 29.81 | 1.21 | 92.55 | 51.92 | 29.57 |
| LLaMA3-8B-Instruct | 44.47 | 25.33 | 32.27 | 2.86 | 80.92 | 51.57 | 42.54 |
| LLaMA3-70B-Instruct | 56.82 | 38.80 | 46.11 | 2.38 | 86.38 | 66.56 | 56.40 |
| Qwen2.5-7B-Instruct | 55.50 | 55.30 | 55.40 | 1.64 | 91.43 | 69.32 | 57.53 |
| Qwen2.5-72B-Instruct | 61.90 | 70.36 | 65.86 | 1.36 | 94.10 | 73.99 | 64.15 |
| DeepSeek-V3 | 56.47 | <u>71.32</u> | 63.03 | <u>1.20</u> | 95.26 | 74.76 | 62.76 |
| <i>Specialized Models</i> | | | | | | | |
| xLAM-7b-fc-r | 0.27 | 0.54 | 0.36 | 0.80 | 88.15 | 11.45 | 4.60 |
| gorilla-openfunctions-v2 | 10.11 | 7.13 | 8.36 | 2.31 | 70.83 | 37.90 | 19.23 |
| ToolLLaMA-2-7b-v2 | 1.89 | 1.34 | 1.57 | 2.29 | 58.77 | 18.29 | 5.01 |
| Mistral-Interact | 4.99 | 4.16 | 4.53 | 1.95 | 83.10 | 25.47 | 9.89 |
| <i>Ours</i> | | | | | | | |
| ASKToACT-LoRA-SFT-7B | 57.68 (↑2.18) | 63.41 (↑8.11) | 60.41 (↑5.01) | 1.48 (↓0.16) | <u>96.05</u> (↑4.62) | <u>81.42</u> (↑12.10) | <u>68.71</u> (↑11.18) |
| ASKToACT-Full-SFT-7B | <u>63.88</u> (↑8.38) | 68.10 (↑12.80) | <u>65.92</u> (↑10.52) | 1.53 (↓0.11) | 97.37 (↑5.94) | 84.55 (↑15.23) | 73.12 (↑15.59) |

Table 3: Main results.

| Metric | Qwen2.5-7B-Instruct | ASKToACT-LoRA-SFT-7B | ASKToACT-Full-SFT-7B |
|-------------------------------|---------------------|----------------------|----------------------|
| Task Completion Rate (%) | 86.67 | 93.33 (↑6.66) | 96.67 (↑10.00) |
| Intent Coverage Rate (%) | 65.85 | 75.61 (↑9.76) | 85.37 (↑19.52) |
| Interaction Rounds | 3.20 | 2.73 (↓0.47) | 2.60 (↓0.60) |
| User Satisfaction Score (1–5) | 3.80 | 4.40 (↑0.60) | 4.61 (↑0.81) |

Table 4: Human evaluation results. All metrics are averaged across participants.

As shown in Table 5, both the LoRA and fully fine-tuned variants of our method demonstrate strong performance. The LoRA variant achieves a CPS of 60.23% and PRS of 64.81%, outperforming all open-source baselines and even surpassing some commercial closed-source models. The fully fine-tuned variant pushes this further, reaching a CPS of 62.96% and PRS of 69.45%, comparable to GPT-4o. These results highlight that our method generalizes effectively to unseen domains without relying on memorization of training data. Instead, it acquires transferable principles for intent clarification and tool invocation.

3.6 Ablation Study

To assess the contribution of each component in our method, we conducted a comprehensive ablation study comparing three model configurations: (1) ASKToACT-LoRA-SFT-7B model, (2) a variant without error-correction augmented dialogue data (i.e., trained only with basic intent clarification data

using the same LoRA configurations), and (3) the untrained base model (Qwen2.5-7B-Instruct). We randomly selected 50 unspecified user queries from the test set and computed the error rates for five error types identified in §2.2.

As shown in Table 6, compared to the untrained base model, the model trained solely on basic intent clarification data significantly reduce all five error types, confirming the effectiveness of clarification training. Incorporating error-correction augmented dialogues and self-correction training yields further improvements. The Clearly Stated Intent Clarification rate and Redundant Clarification rate both decrease from 9.09% to 6.80%, suggesting that the model becomes more effective at avoiding unnecessary clarification. While Imprecise and Irrelevant Clarification rates show slight increases, likely due to additional interaction turns introduced by self-correction attempts, this trade-off is justified by the substantial reduction in Incomplete Clarification rate (from 38.00% to 32.00%), which is critical for

| LLM | Intent Clarification Quality | | | | Tool Invocation Accuracy | | |
|---------------------------|------------------------------|-----------------------|-----------------------|--------------|--------------------------|-----------------------|----------------------|
| | ICR↑ | CE↑ | CPS↑ | IR↓ | SCR↑ | TSS↑ | PRS↑ |
| <i>Closed-Source LLMs</i> | | | | | | | |
| Claude3.5-Haiku | 61.07 | 29.88 | 40.13 | 4.20 | 73.68 | 66.15 | 46.59 |
| Claude3.5-Sonnet | 69.74 | 38.10 | 49.28 | 3.29 | 84.96 | 76.05 | 54.09 |
| GPT-3.5 | 44.19 | 44.42 | 44.30 | 1.72 | 98.27 | 89.28 | 45.50 |
| GPT-4 | 63.60 | 44.73 | 52.52 | 2.57 | 93.90 | 90.52 | 63.26 |
| GPT-4o-mini | <u>70.86</u> | 49.60 | 58.35 | 2.63 | 95.22 | 89.63 | 69.44 |
| GPT-4o | 72.41 | 53.37 | <u>61.45</u> | 2.27 | 96.24 | 92.22 | 69.56 |
| <i>Open-Source LLMs</i> | | | | | | | |
| Mistral-7B-Instruct-v0.3 | 55.13 | 28.94 | 37.96 | 3.15 | 77.34 | 68.73 | 49.35 |
| LLaMA3-8B-Instruct | 62.81 | 29.99 | 40.59 | 3.66 | 78.86 | 69.22 | 45.52 |
| LLaMA3-70B-Instruct | 67.14 | 35.34 | 46.30 | 3.44 | 84.76 | 79.08 | 55.53 |
| Qwen2.5-7B-Instruct | 64.79 | 38.43 | 48.25 | 3.00 | 92.99 | 86.19 | 61.86 |
| Qwen2.5-72B-Instruct | 68.64 | 43.87 | 53.53 | 2.83 | 92.48 | 90.25 | 63.85 |
| DeepSeek-V3 | 60.11 | 42.24 | 49.62 | 2.55 | 92.17 | 83.10 | 58.33 |
| <i>Specialized Models</i> | | | | | | | |
| xLAM-7b-fc-r | 0.34 | 0.56 | 0.42 | <u>2.08</u> | 91.46 | 14.29 | 5.43 |
| gorilla-openfunctions-v2 | 44.53 | 22.27 | 29.69 | 3.41 | 69.92 | 52.36 | 21.67 |
| ToolLLaMA-2-7b-v2 | 2.76 | 2.25 | 2.48 | 2.19 | 98.98 | 42.65 | 2.07 |
| Mistral-Interact | 35.38 | 15.63 | 21.69 | 4.27 | 64.43 | 18.60 | 2.94 |
| <i>Ours</i> | | | | | | | |
| ASKToACT-LoRA-SFT-7B | 68.87 (↑4.08) | <u>53.52</u> (↑15.09) | 60.23 (↑11.98) | 2.82 (↑0.18) | <u>99.59</u> (↑6.63) | 96.44 (↑10.25) | 64.81 (↑2.95) |
| ASKToACT-Full-SFT-7B | 69.90 (↑5.11) | 57.27 (↑18.84) | 62.96 (↑14.71) | 2.72 (↑0.28) | 99.70 (↑6.64) | <u>96.41</u> (↑10.22) | <u>69.45</u> (↑7.59) |

Table 5: OOD generalization performance comparison.

| Method | Error Rate (%) | | | | | CPS↑ | PRS↑ |
|--|--|----------------------------|-----------------------------|----------------------------|-----------------------------|-------|-------|
| | Clearly Stated Intent Clarification | Imprecise Clarification | Irrelevant Clarification | Redundant Clarification | Incomplete Clarification | | |
| ASKToACT-LoRA-SFT-7B | 6.80 | 11.65 | 8.74 | 6.80 | 32.00 | 61.51 | 69.50 |
| w/o Error-Correction Augmented Dialogue Data | 9.09 | 6.49 | 6.49 | 9.09 | 38.00 | 58.93 | 66.90 |
| w/o Training (Base Model) | 12.43 | 12.37 | 9.29 | 11.34 | 44.00 | 53.00 | 59.18 |

Table 6: Ablation study. The first four error types calculated as the proportion of interaction turns containing specific errors among all interaction turns, while the last error type measures the proportion of queries in which not all unspecified intents are successfully clarified.

enabling accurate tool invocation. These improvements in clarification behavior are further reflected in downstream performance. The PRS increases from 59.18% to 66.90%, and finally to 69.50%, indicating that the enhanced clarification quality translates into more accurate tool invocation.

4 Related Work

Our work relates to three areas: *tool learning*, *user intent clarification*, and *self-correction*. Tool learning equips LLMs with external capabilities but typically assumes explicit user intents. Intent clarification addresses ambiguous queries, yet existing datasets often rely on manual annotation. Self-correction has shown promise in mathematical reasoning but remains underexplored for intent understanding. We unify these directions through a self-correcting clarification framework. A full review of related work is provided in Appendix A.

5 Conclusion

In this work, we presented **ASKToACT**, a self-correcting clarification framework for tool learning that addresses the critical challenges of data scalability and error handling in clarification interactions. Our key contribution lies in leveraging the inherent structure of tool learning datasets to enable automated construction of high-quality intent clarification data, while introducing a novel self-correction mechanism for robust clarification. Experimental results demonstrate that our method not only achieves superior performance in intent clarification and tool invocation but also exhibits strong generalization to unseen APIs and diverse model architectures. We hope that our work will provide valuable insights for developing more effective and reliable intent clarification mechanisms in human-LLM interaction systems.

Limitations

While our work demonstrates promising results in handling unspecified queries, several important limitations warrant discussion:

Dataset and Training Our method heavily relies on existing tool learning datasets, which may not fully capture the diversity and complexity of real-world user intents. The parameter removal process, although systematic, might not perfectly simulate natural query ambiguity patterns. Additionally, our current approach to error-correction augmentation focuses on pre-defined error types, potentially missing other important error patterns that emerge in real-world interactions.

Interaction Dynamics We have not yet explored scenarios where intents must be inferred from previous tool invocation results, limiting our framework’s ability to handle context-dependent queries.

Evaluation Limitations While our multi-level evaluation framework provides comprehensive assessment, it may not fully capture the complexity of real-world deployment scenarios, particularly in terms of user patience, time constraints, and varying expertise levels. The current metrics might not sufficiently measure the user experience aspects of the clarification process.

Ethics Statement

We acknowledge that all authors are informed about and adhere to the ACL Code of Ethics and the Code of Conduct.

Use of AI-Generated Content In our research, we utilize LLMs to generate intent clarification dialogues based on existing tool learning datasets. All AI-generated content has been thoroughly verified by the authors to ensure quality and appropriateness. We have implemented rigorous quality control mechanisms to filter out inappropriate or low-quality generations. The paper clearly discloses all instances where AI systems contributed to content generation.

Data Sources The tool learning datasets used in our experiments are derived from publicly available sources, including open-source repositories and publicly released benchmarks. We have made reasonable efforts to ensure that these data sources do not contain personally identifiable information or legally protected content. However, we cannot

guarantee that they are entirely free from socially harmful or biased language. Any potential biases in the original datasets may propagate to our results.

Broader Impact Our work aims to enhance models’ ability to handle ambiguous user queries in tool-use scenarios. This may extend the applicability of AI systems to a wider range of real-world scenarios. However, such improvements in intent clarification and tool-use capabilities could also enable models to act with limited human oversight, posing both opportunities and risks depending on the deployment context.

References

- Chinmaya Andukuri, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah D Goodman. 2024. Star-gate: Teaching language models to ask clarifying questions. *arXiv preprint arXiv:2403.19154*.
- Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew White, and Philippe Schwaller. 2023. [Augmenting large language models with chemistry tools](#). In *NeurIPS 2023 AI for Science Workshop*.
- Hongqiao Chen, Kexun Zhang, Lei Li, and William Yang Wang. 2023. [Tooldec: Syntax error-free and generalizable tool use for LLMs via finite-state decoding](#). In *The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS’23*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Yu Gu, Yiheng Shu, Hao Yu, Xiao Liu, Yuxiao Dong, Jie Tang, Jayanth Srinivasa, Hugo Latapie, and Yu Su. 2024. [Middleware for LLMs: Tools are instrumental for language agents in complex environments](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7646–7663, Miami, Florida, USA. Association for Computational Linguistics.
- Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. 2023. [Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 45870–45894. Curran Associates, Inc.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). *Preprint*, arXiv:2106.09685.
- Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023. [Large](#)

| | | |
|-----|---|-----|
| 641 | language models can self-improve. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1051–1068, Singapore. Association for Computational Linguistics. | 699 |
| 642 | | 700 |
| 643 | | 701 |
| 644 | | 702 |
| 645 | Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2024a. Large language models cannot self-correct reasoning yet. In <i>The Twelfth International Conference on Learning Representations</i> . | 703 |
| 646 | | 704 |
| 647 | | 705 |
| 648 | | 706 |
| 649 | | 707 |
| 650 | | |
| 651 | Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao Wan, Neil Zhenqiang Gong, and Lichao Sun. 2024b. Meta-tool benchmark for large language models: Deciding whether to use tools and which to use. In <i>The Twelfth International Conference on Learning Representations</i> . | 708 |
| 652 | | 709 |
| 653 | | 710 |
| 654 | | 711 |
| 655 | | 712 |
| 656 | | 713 |
| 657 | | 714 |
| 658 | Tatsuro Inaba, Hirokazu Kiyomaru, Fei Cheng, and Sadao Kurohashi. 2023. MultiTool-CoT: GPT-3 can use multiple external tools with chain of thought prompting. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1522–1532, Toronto, Canada. Association for Computational Linguistics. | 715 |
| 659 | | 716 |
| 660 | | |
| 661 | | 717 |
| 662 | | 718 |
| 663 | | 719 |
| 664 | | 720 |
| 665 | | 721 |
| 666 | | 722 |
| 667 | | 723 |
| 668 | | 724 |
| 669 | | 725 |
| 670 | | |
| 671 | 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 <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 3102–3116, Singapore. Association for Computational Linguistics. | 726 |
| 672 | | 727 |
| 673 | | 728 |
| 674 | | 729 |
| 675 | | 730 |
| 676 | | 731 |
| 677 | | 732 |
| 678 | | 733 |
| 679 | Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, Rithesh Murthy, Liangwei Yang, Silvio Savarese, Juan Carlos Niebles, Huan Wang, Shelby Heinecke, and Caiming Xiong. 2024. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets. In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 54463–54482. Curran Associates, Inc. | 734 |
| 680 | | 735 |
| 681 | | 736 |
| 682 | | 737 |
| 683 | | 738 |
| 684 | | |
| 685 | | 739 |
| 686 | | 740 |
| 687 | | 741 |
| 688 | | 742 |
| 689 | | 743 |
| 690 | | 744 |
| 691 | | 745 |
| 692 | | |
| 693 | | 746 |
| 694 | | 747 |
| 695 | | 748 |
| 696 | | 749 |
| 697 | | 750 |
| 698 | | |
| | Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. <i>arXiv preprint arXiv:2112.09332</i> . | 751 |
| | | 752 |
| | | 753 |
| | | 754 |
| | | 755 |
| | Kangyun Ning, Yisong Su, Xueqiang Lv, Yuanzhe Zhang, Jian Liu, Kang Liu, and Jinan Xu. 2024. Wtu-eval: A whether-or-not tool usage evaluation benchmark for large language models. <i>Preprint</i> , arXiv:2407.12823. | |
| | | |
| | Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. 2023. Art: Automatic multi-step reasoning and tool-use for large language models. <i>Preprint</i> , arXiv:2303.09014. | |
| | | |
| | 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> . | |
| | | |
| | Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Zhong Zhang, Jie Zhou, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. Tell me more! towards implicit user intention understanding of language model driven agents. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1088–1113, Bangkok, Thailand. Association for Computational Linguistics. | |
| | | |
| | Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. 2024. ToolLLM: Facilitating large language models to master 16000+ real-world APIs. In <i>The Twelfth International Conference on Learning Representations</i> . | |
| | | |
| | Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics. | |
| | | |
| | Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 68539–68551. Curran Associates, Inc. | |
| | | |
| | Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugging-GPT: Solving AI tasks with chatGPT and its friends in hugging face. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> . | |
| | | |
| | Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. 2024. Taskbench: Benchmarking large language models for task automation. In <i>Advances in Neural Information Processing Systems</i> , | |

| | | |
|-----|---|-----|
| 756 | volume 37, pages 4540–4574. Curran Associates, Inc. | 813 |
| 757 | | 814 |
| 758 | Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023. Replug: Retrieval-augmented black-box language models . <i>Preprint</i> , arXiv:2301.12652. | 815 |
| 759 | | 816 |
| 760 | | 817 |
| 761 | | 818 |
| 762 | | 819 |
| 763 | Jeonghoon Shim, Gyuhyeon Seo, Cheongsu Lim, and Yohan Jo. 2025. Tooldial: Multi-turn dialogue generation method for tool-augmented language models . In <i>The Thirteenth International Conference on Learning Representations</i> . | 820 |
| 764 | | 821 |
| 765 | | 822 |
| 766 | | 823 |
| 767 | | |
| 768 | Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, Ye Tian, and Sujian Li. 2023. Restgpt: Connecting large language models with real-world restful apis . <i>Preprint</i> , arXiv:2306.06624. | 824 |
| 769 | | 825 |
| 770 | | 826 |
| 771 | | 827 |
| 772 | | |
| 773 | Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. 2023. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases . <i>Preprint</i> , arXiv:2306.05301. | 828 |
| 774 | | 829 |
| 775 | | 830 |
| 776 | | |
| 777 | Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguerre-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. Lamda: Language models for dialog applications . <i>Preprint</i> , arXiv:2201.08239. | 831 |
| 778 | | 832 |
| 779 | | 833 |
| 780 | | 834 |
| 781 | | 835 |
| 782 | | 836 |
| 783 | | 837 |
| 784 | | 838 |
| 785 | | |
| 786 | | 839 |
| 787 | | 840 |
| 788 | | 841 |
| 789 | | 842 |
| 790 | | 843 |
| 791 | | |
| 792 | | |
| 793 | | |
| 794 | | |
| 795 | | |
| 796 | | |
| 797 | | |
| 798 | Hongru Wang, Yujia Qin, Yankai Lin, Jeff Z Pan, and Kam-Fai Wong. 2024a. Empowering large language models: Tool learning for real-world interaction. In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2983–2986. | |
| 799 | | |
| 800 | | |
| 801 | | |
| 802 | | |
| 803 | | |
| 804 | Wenxuan Wang, Juluan Shi, Chaozheng Wang, Cheryl Lee, Youliang Yuan, Jen-tse Huang, and Michael R Lyu. 2024b. Learning to ask: When llms meet unclear instruction. <i>arXiv preprint arXiv:2409.00557</i> . | |
| 805 | | |
| 806 | | |
| 807 | | |
| 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>Advances in Neural Information Processing Systems</i> , | |
| 809 | | |
| 810 | | |
| 811 | | |
| 812 | | |
| | Yuchen Yan, Jin Jiang, Yang Liu, Yixin Cao, Xin Xu, Mengdi zhang, Xunliang Cai, and Jian Shao. 2024. S³c-math: Spontaneous step-level self-correction makes large language models better mathematical reasoners . <i>Preprint</i> , arXiv:2409.01524. | |
| | | |
| | Che Zhang, Zhenyang Xiao, Chengcheng Han, Yixin Lian, and Yuejian Fang. 2024. Learning to check: Unleashing potentials for self-correction in large language models . <i>Preprint</i> , arXiv:2402.13035. | |
| | | |
| | Kechi Zhang, Huangzhao Zhang, Ge Li, Jia Li, Zhuo Li, and Zhi Jin. 2023. Toolcoder: Teach code generation models to use api search tools . <i>Preprint</i> , arXiv:2305.04032. | |
| | | |
| | Michael JQ Zhang and Eunsol Choi. 2023. Clarify when necessary: Resolving ambiguity through interaction with lms. <i>arXiv preprint arXiv:2311.09469</i> . | |
| | | |
| | Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , Bangkok, Thailand. Association for Computational Linguistics. | |
| | | |
| | Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. Toolqa: A dataset for llm question answering with external tools . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 50117–50143. Curran Associates, Inc. | |
| | | |

Appendix

A Related Work

A.1 Tool Learning

Tool learning can effectively alleviate the inherent limitations of LLMs through dynamic interaction with external tools (Schick et al., 2023; Tang et al., 2023; Shen et al., 2023; Qin et al., 2024; Wang et al., 2024a). While LLMs acquire limited knowledge during the pre-training phase, tools such as integrated search engines (Nakano et al., 2021; Komeili et al., 2022; Schick et al., 2023; Zhang et al., 2023; Shi et al., 2023; Paranjape et al., 2023) and databases (Thoppilan et al., 2022; Patil et al., 2023; Hao et al., 2023; Zhuang et al., 2023; Chen et al., 2023; Gu et al., 2024) enable real-time access to up-to-date information beyond the training data. In addition, LLMs often struggle with complex mathematical operations, code generation, and domain-specific tasks (Inaba et al., 2023; Bran et al., 2023), which can be enhanced through dedicated tools.

Existing evaluation benchmarks for reliable tool usage (Huang et al., 2024b; Patil et al., 2023; Ning et al., 2024) focus on explicit and unambiguous user queries, leaving the challenges of handling unspecified intents in real-world scenarios largely unexplored.

A.2 User Intent Clarification

When interacting with users, understanding user intents is crucial, especially when intents are implicit or unspecified. Zhang and Choi (2023) shows that unspecified user intents in queries should be clarified through interaction. The STaR-GATE framework (Andukuri et al., 2024) introduces a systematic approach to question formulation by simulating diverse clarification scenarios. Qian et al. (2024) applied several strategies in conversation record construction and leveraged the generated data to fine-tune the model, enhancing the ability to formulate targeted questions.

However, the construction of high-quality datasets for training and evaluation still remains challenging. Qian et al. (2024) constructed a benchmark for daily scenarios, while Wang et al. (2024b) focuses on tool learning scenarios, but they both relied on manual annotation. Our work introduces an automated pipeline for dataset construction, enabling better scalability.

A.3 Self-Correction Mechanism

Early work on self-correction (Huang et al., 2023; Madaan et al., 2023) primarily focused on post-correction, using feedback to improve model outputs after they are generated. However, Huang et al. (2024a) found that in the absence of standardized answers, such post-correction has limited effect. This finding prompted a shift in research focus to real-time self-correction, i.e., dynamically identifying and correcting errors during the reasoning process.

Self-correction has achieved success in mathematical reasoning, where Yan et al. (2024) and Zhang et al. (2024) introduce step-level and multi-granular correction strategies. We extend these approaches to user intent clarification, enabling real-time correction during the clarification process.

B Datasets

Our dataset is constructed based on two existing tool learning datasets: xlam-function-calling-60k and TaskBench. We describe their characteristics below.

xlam-function-calling-60k This dataset comprises functionally executable APIs extracted from Python libraries and RESTful services, rather than being manually defined. The APIs span 21 functional categories, covering a broad range of domains such as information retrieval, and computational tools. In total, the dataset contains 3,673 APIs and 60,000 samples. An example is shown below:

```
{
  "query": "List titles originally aired on
networks '1' and '8', released after 2010,
sorted by release date in descending order.",
  "tools": [
    {
      "name": "list_titles",
      "description": "Fetches a listing of
titles that match specified parameters from the
Watchmode API.",
      "parameters": {
        "genres": {
          "description": "Filter
results to only include certain genre(s). Pass
in a single genre ID or multiple comma-separated
IDs. Default is '4,9'.",
          "type": "str",
          "default": "4,9"
        },
        "limit": {
          "description": "Set how many
titles to return per page. Default and maximum
is 250.",
```



```

947         "type": "int",
948         "default": "250"
949     },
950     "source_ids": {
951         "description": "Filter the
952 results to titles available on specific sources
953 by passing individual IDs or multiple comma-
954 separated IDs. Default is '23,206'. Note: Only a
955 single region can be set if this is populated.",
956         "type": "str",
957         "default": "23,206"
958     },
959     "source_types": {
960         "description": "Filter
961 results to only include titles available on
962 specific types of sources (e.g., subscription,
963 free). Default is 'sub,free'. Note: Only a
964 single region can be set if this is populated.",
965         "type": "str",
966         "default": "sub,free"
967     },
968     "types": {
969         "description": "Filter
970 results to only include titles available on
971 specific types of sources (e.g., subscription,
972 free). Default is 'sub,free'. Note: Only a
973 single region can be set if this is populated.",
974         "type": "str",
975         "default": "movie,tv_series"
976     },
977     "regions": {
978         "description": "Filter
979 results to only include sources active in
980 specific regions. Currently supported regions:
981 US, GB, CA, AU. Default is 'US'. Note: Only a
982 single region can be set if source_ids or source
983 _types are populated.",
984         "type": "str",
985         "default": "US"
986     },
987     "sort_by": {
988         "description": "Sort order
989 of results. Possible values include: relevance_
990 desc, relevance_asc, popularity_desc, popularity
991 _asc, release_date_desc, release_date_asc, title
992 _desc, title_asc. Default is 'relevance_desc'.",
993         "type": "str",
994         "default": "relevance_desc"
995     },
996     "page": {
997         "description": "Set the page
998 of results to return. Default is 1.",
999         "type": "int",
1000         "default": "1"
1001     },
1002     "network_ids": {
1003         "description": "Filter
1004 results to titles that originally aired on
1005 specific TV networks by passing individual IDs
1006 or multiple comma-separated IDs. Default is
1007 '1,8,12'.",
1008         "type": "str",
1009         "default": "1,8,12"
1010     },
1011     "release_date_start": {
1012         "description": "Filter
1013 results to only include titles released on or
1014 after a specific date. Format: YYYYMMDD. Default
1015 is 20010101.",

```

```

        "type": "int",
        "default": "20010101"
    },
    "release_date_end": {
        "description": "Filter
results to only include titles released on or
before a specific date. Format: YYYYMMDD.
Default is 20201211.",
        "type": "int",
        "default": "20201211"
    }
},
],
"answers": [
    {
        "task": "list_titles",
        "arguments": [
            { "name": "network_ids", "value":
"1,8" },
            { "name": "release_date_start",
"value": 20110101 },
            { "name": "sort_by", "value": "
release_date_desc" }
        ]
    }
]
}

```

Taskbench In contrast, TaskBench defines three tool-use scenarios across distinct application domains: (1) Hugging Face tools, simulating a subset of model functionalities available on the Hugging Face platform—such as summarization, translation, and classification, with 23 APIs and 12,217 samples; (2) Multimedia tools, simulating representative functionalities of multimodal systems—such as video editing and image transformation, with 40 APIs and 8,904 samples; (3) Daily Life APIs, simulating everyday user-facing applications—such as ticket booking, food ordering and schedule management, with 40 APIs and 7,150 samples. All APIs in TaskBench are manually constructed. Representative examples from each domain are shown below:

Hugging Face tools

```

{
    "query": "I'm currently analyzing a
particular text, 'John works at Google in
Mountain View, California.' Can you assist me in
identifying the named entities and marking the
part-of-speech tags within this text?",
    "tools": [
        {
            "id": "Token Classification",
            "desc": "Token classification is a
natural language understanding task in which a
label is assigned to some tokens in a text. Some
popular token classification subtasks are Named
Entity Recognition (NER) and Part-of-Speech (
PoS) tagging. NER models could be trained to
identify specific entities in a text, such as

```

| | | | |
|------|--|--|------|
| 1079 | dates, individuals and places; and PoS tagging | "abc.txt" | 1147 |
| 1080 | would identify, for example, which words in a |] | 1148 |
| 1081 | text are verbs, nouns, and punctuation marks.", | } | 1149 |
| 1082 | "input-type": [|] | 1150 |
| 1083 | "text or text file" | } | 1151 |
| 1084 |], | | |
| 1085 | "output-type": [| | 1152 |
| 1086 | "text or text file" | | 1153 |
| 1087 |] | | |
| 1088 | } | | 1154 |
| 1089 |], | | 1155 |
| 1090 | "answers": [| "query": "I have a busy day ahead. Could you | 1156 |
| 1091 | { | assist me by logging into an online meeting | 1157 |
| 1092 | "task": "Token Classification", | regarding 'Smart Home Devices'? After the | 1158 |
| 1093 | "arguments": [| meeting, can you facilitate a video call with my | 1159 |
| 1094 | "'John works at Google in | friend on +1234567666?", | 1160 |
| 1095 | Mountain View, California.'" | "tools": [| 1161 |
| 1096 |] | { | 1162 |
| 1097 | } | "id": "attend_meeting_online", | 1163 |
| 1098 |] | "desc": "Attend a meeting online | 1164 |
| 1099 | } | about a specific topic", | 1165 |
| | | "parameters": [| 1166 |
| | | { | 1167 |
| | | "name": "topic", | 1168 |
| | | "type": "string", | 1169 |
| | | "desc": "The topic of the | 1170 |
| 1100 | | meeting" | 1171 |
| 1101 | <i>Multimedia tools</i> | } | 1172 |
| 1102 | { |] | 1173 |
| 1103 | "query": "I've recently conducted an | }, | 1174 |
| 1104 | interview and have recorded it in 'interview.wav | { | 1175 |
| 1105 | ' audio file. Can you assist me in transcribing | "id": "make_video_call", | 1176 |
| 1106 | it to a text document, so I can refer to it | "desc": "Make a video call to a | 1177 |
| 1107 | easily in the future? Besides, I'm dealing with | specific phone number", | 1178 |
| 1108 | an article titled 'abc.txt' and I want to have a | "parameters": [| 1179 |
| 1109 | fresh iteration of this text so that it will be | { | 1180 |
| 1110 | unique. Would you be able to employ the Article | "name": "phone_number", | 1181 |
| 1111 | Spinner tool to facilitate this?", | "type": "string", | 1182 |
| 1112 | "tools": [| "desc": "The phone number to | 1183 |
| 1113 | { | make the video call to" | 1184 |
| 1114 | "id": "Audio-to-Text", | } | 1185 |
| 1115 | "desc": "Transcribes speech from an |] | 1186 |
| 1116 | audio file into text.", | } | 1187 |
| 1117 | "input-type": [|], | 1188 |
| 1118 | "audio or audio file" | "answers": [| 1189 |
| 1119 |], | { | 1190 |
| 1120 | "output-type": [| "task": "attend_meeting_online", | 1191 |
| 1121 | "text or text file" | "arguments": [| 1192 |
| 1122 |] | { "name": "topic", "value": " | 1193 |
| 1123 | }, | Smart Home Devices" } | 1194 |
| 1124 | { |] | 1195 |
| 1125 | "id": "Article Spinner", | { | 1196 |
| 1126 | "desc": "Rewrites a given article | "task": "make_video_call", | 1197 |
| 1127 | using synonyms and syntax changes to create a | "arguments": [| 1198 |
| 1128 | new, unique version.", | { "name": "phone_number", "value | 1199 |
| 1129 | "input-type": [| ": "+1234567666" } | 1200 |
| 1130 | "text or text file" |] | 1201 |
| 1131 |], | } | 1202 |
| 1132 | "output-type": [|] | 1203 |
| 1133 | "text or text file" | } | 1204 |
| 1134 |] | | |
| 1135 | } | | |
| 1136 |], | | |
| 1137 | "answers": [| | 1205 |
| 1138 | { | We verified that there is no overlap between | 1206 |
| 1139 | "task": "Audio-to-Text", | APIs in xlam-function-calling-60k and those in | 1207 |
| 1140 | "arguments": [| TaskBench. | |
| 1141 | "interview.wav" | | |
| 1142 |] | Licensing Both datasets are publicly accessi- | 1208 |
| 1143 | }, | ble: xlam-function-calling-60k follows the Cre- | 1209 |
| 1144 | { | ative Commons Attribution 4.0 License (CC BY) | 1210 |
| 1145 | "task": "Article Spinner", | while TaskBench is released under the Apache 2.0 | 1211 |
| 1146 | "arguments": [| | |

License. We comply with their respective licenses in using and extending the data.

C Intent Clarification Dataset Curation

C.1 Unspecified Query Generation

C.1.1 Prompt for Unspecified Query Generation

Given an original query and its tool invocation solution, the following prompt guides GPT-4o to generate unspecified queries by analyzing parameters and systematically removing them.

System Prompt

You are a query transformation assistant. Your task is to modify the original user query by removing or abstracting specific parameters marked with ``removed``: "true", while maintaining the overall structure and clarity of the original query. The resulting query (``unspecified_query``) should reflect the general intent of the user but omit or obscure the specific details of the removed parameters.

Input:

- ``original_query``: The complete textual description of the user's original request.
- ``answers``: A detailed record specifying the APIs and parameters required to fulfill the original query. Each parameter in this record includes:
 - ``removed``: A boolean ("true" or "false") indicating whether this parameter should be removed or abstracted during the transformation process.
 - Other relevant metadata, such as the parameter's value.
- ``tools``: Documentation or descriptive details about the tools referenced in the query, including their parameters and usage instructions.

Transformation Rules:

- Identify the parameters to be removed or abstracted:
 - Focus on parameters where ``removed``: "true".
 - Identify the full range of corresponding expressions in the query, ensuring all references to the parameter are appropriately handled.
- Apply the appropriate transformation strategy to each parameter marked as ``removed``: "true":
 - Complete Removal**: The parameter is entirely removed when it has no significant impact on the remaining content of the query. However, this should not be used if the parameter is optional. Also, if the same tool is called multiple times, the parameters should not be removed. Instead, they should be abstracted.
 - Semantic Abstraction**: If the parameter influences the meaning or structure of the query, replace its value with a more general or abstract term.

- Partial Obfuscation**: If the elements of a matrix or list are presented separately in the query (e.g., discrete values like quantities or items) and need to be constructed or inferred, only one element from the matrix or list should be removed or abstracted. The remaining elements should stay intact. This can still be done using either complete removal or semantic abstraction, while leaving other relevant elements unchanged.
- Ensure textual and structural coherence:
 - After transformation, ensure that the ``unspecified_query`` remains readable, logically consistent, and grammatically correct.
- Avoid explicitly stating "unspecified" or "unknown" values:
 - Do not use terms like "unspecified", "unknown", or "ambiguous" in the ``unspecified_query``.
 - Instead, naturally omit or generalize the missing details without drawing attention to their absence.
- Retain the rest of the query:
 - Leave unchanged the parts of the query that are not marked for removal, maintaining consistency in format and information.

Output:

Return a JSON object containing:

- ``unspecified_query``: The transformed query string with removed/abstracted parameters.
- ``key_info``: A JSON array (or object) documenting all parameters, containing the following fields:
 - ``original_value``: The expression of the parameter as it appears in the ``original_query`` (not the value in the ``answers``).
 - ``current_value``: The transformed value of the parameter in the ``unspecified_query``.
 - ``removed``: The boolean flag indicating whether the parameter was removed.

C.1.2 Transformation Record Format

For each unspecified query, we maintain a transformation record in the following JSON structure:

```
{
  "original_query": string,
  "unspecified_query": string,
  "key_info": [
    {
      "name": [API_name],
      "arguments": {
        [param_name]: {
          "original_value": string,
          "current_value": string,
          "removed": boolean
        },
        ...
      },
      ...
    },
    ...
  ]
}
```

C.1.3 Human Verification

To further ensure the quality of generated unspecified queries, we perform human verification on

| Metric | Naturalness↑ | Consistency↑ | Necessity↑ | Complexity↑ | Diversity↑ | Acceptance Rate (%)↑ |
|--------|--------------|--------------|------------|-------------|------------|----------------------|
| Score | 4.61/5 | 4.80/5 | 4.03/5 | 3.87/5 | 4.54/5 | 95.92 |

Table 7: Human verification results of unspecified query generation. The first five metrics are rated on a 1–5 scale, while Acceptance Rate is reported as a percentage. All metrics are averaged across participants.

400 randomly sampled queries. Three graduate students with NLP backgrounds independently assessed each query based on six criteria: Naturalness (fluency and linguistic coherence), Consistency (uniformity of transformation), Necessity (need for clarification), Complexity (difficulty of clarification), Diversity (range of parameter types and domains), and Acceptance Rate (overall acceptability). Results are shown in Table 7.

C.2 Clarification Dialogue Construction

We divide the clarification dialogue construction process into two steps: GPT-4o-dependent content generation and template-based dialogue assembly. The essential content is generated using GPT-4o and stored a structured transformation record (see Appendix C.1.2 for format details). The information encoded in this record is then used to deterministically assemble the final dialogue through predefined templates, without further reliance on GPT-4o.

C.2.1 GPT-4o-Dependent Content Generation

Task Decomposition We employ the following prompt to guide GPT-4o in decomposing user queries into subtasks:

System Prompt

You are a smart task decomposition assistant. Your goal is to break down the user's main task into smaller, manageable subtasks. Please follow the instructions below.

You will receive a JSON-formatted input containing:

- ``query``: A description of the main task the user wants to accomplish.
- ``tools``: A list of APIs available to solve the task, each with a unique identifier and a description of its functionality. Note: The APIs are provided in the exact order necessary to resolve the task.

Task Decomposition:

1. Analyze the query to understand the user's main task.
2. Break it down into smaller, manageable subtasks that can be handled using the provided APIs. Ensure that each subtask is completed by calling one of the APIs in the exact order they are listed.

Your output should be a JSON object with the following structure:

```
{
  "tool_steps": [
    "Step <number>: <subtask description>
    using <API name>.",
    ...
  ]
}
```

The decomposition result is added to the transformation record as a new field `"tool_steps"`.

Clarification Question Generation We use the following prompt to generate clarification questions for unspecified parameters:

System Prompt

You are an assistant responsible for generating clarification questions for missing information in the user's query.

Input:

The input should contain the following fields:

- ``original_query``: A complete user task description.

- ``unspecified_query``: A user task description missing some key information.
- ``tools``: Documentation or descriptive details about the tools referenced in the query, including their parameters and usage instructions.
- ``key_info``: This should record the APIs and parameters needed to solve the user task, including information about any missing parameters.

- ``original_value``: The original value of the parameter in the ``original_query``.
- ``current_value``: The current value of the parameter in the ``unspecified_query``.
- ``removed``: Indicates if the parameter's value is clear ("false") or unspecified ("true").

Task Requirements:

For each parameter where the field ``removed`` is set to true, you are to generate a clarification question.

- If multiple API calls rely on the same missing information, form a single combined question to efficiently gather the required details, rather than asking multiple separate questions.
- Each question should focus on gathering one specific piece of information to improve the precision of the query and avoid ambiguity.
- Do not ask about information that can be inferred from context or API interactions. Only generate clarification questions for details that cannot be deduced from the given context or

API responses.

- Add a `question` field to the corresponding parameter in `key_info`, which contains the generated clarification question.
- Do not modify the `original_query`, `current_value` or any other fields in `key_info`.

Output:
Only output the modified `key_info` in JSON format, ensuring that the question field contains the clarification question for each missing parameter.

The newly generated "key_info" field replaces the original one in the transformation record.

C.2.2 Template-based Dialogue Assembly

Based on the completed transformation record, we automatically construct the dialogue through pre-defined templates.

Task Decomposition We concatenate steps from "tool_steps" to form a comprehensive task decomposition analysis.

Parameter Evaluation For each parameter in the transformation record, we generate evaluation statements using templates based on their removed status:

- For parameters clearly stated in the query, we generate the evaluation that "The parameter [param_name] for API [API_name] has a value of [value]".
- For parameters removed in the query, we generate the evaluation that "The parameter [param_name] for API [API_name] lacks a clear value".

Clarification Interaction Following the API call order, for each removed parameter, we generate a three-part clarification interaction:

- Assistant → User: Ask the clarification question.
- User → Assistant: Provide the original parameter value using templates from Table 8.
- Assistant: Confirm with "Now I know that the parameter [param_name] for API [API_name] has a value of [value]".

Tool Invocation We construct the final tool invocation solution using the "key_info" field from the transformation record, which specifies the sequence of API calls and their associated parameters. The final output is serialized into the following format:

```
[
  {
    "task": [API_name],
    "arguments": [
      {
        "name": [param_name],
        "value": string | number | boolean
      },
      ...
    ]
  },
  ...
]
```

Final Assembly We assemble the complete assistant-user dialogue by sequentially integrating the natural language outputs generated in the previous steps. We begin with a user message that presents the task description and relevant APIs. The assistant's response is then constructed by combining the task decomposition and parameter evaluation. For each missing parameter, we insert a three-part clarification interaction comprising the assistant's question, the user's response, and the assistant's confirmation. This process is repeated until all missing parameters have been clarified. The dialogue concludes with the assistant presenting the full tool invocation solution.

D Self-correction Training

D.1 Error Generation

In our template-based dialogue assembly process (Appendix C.2.2), the sequence of APIs and their required parameters, as recorded in the "key_info" field of the transformation record, implicitly defines the structure of the final dialogue. This insight motivates our error generation strategy. For each selected error type, we first identify a parameter position in the transformation record where the error will be introduced. We then generate the corresponding erroneous behavior and annotate the selected parameter with an "error" field to indicate its error type.

We now describe the generation strategies for each of the five error types in detail.

Clearly Stated Intent Clarification The prompt for generating instances of questioning clearly stated intent is designed as follows:

System Prompt

You are a smart assistant. Your task is to generate a JSON object based on the given input. Please follow these instructions:

Input:

| Tone | Template |
|------------|--|
| Neutral | [value]. The answer is: [value]. Ah, the answer is simply [value]. |
| Friendly | Sure! The answer is [value]. Let me know if you have more questions! I'm glad to help! The answer is absolutely [value]! |
| Dismissive | Honestly, I don't see why this is a big deal, but the answer is [value]. Okay, the answer is: [value]. Hope that helps, I guess. Whatever. The answer is [value]. Not that it matters. |
| Irritated | Listen, the answer is [value]. Just deal with it! Ugh, seriously? The answer is [value]. Can we move on already? Honestly, do you really need me to repeat this? The answer is [value]. I can't believe we're still discussing this! It's infuriating! Enough already! The answer is [value]. Can we please get to the point? I'm tired of this nonsense! It's frustrating! Let's just move on! |

Table 8: Response templates for user with varying tones.

The input should contain the following fields:

- ``original_query``: A complete user task description.
- ``unspecified_query``: A user task description missing some key information.
- ``tools``: Documentation or descriptive details about the tools referenced in the query, including their parameters and usage instructions.
- ``key_info``: This should record the APIs and parameters needed to solve the user task, including information about any missing parameters.

Key Requirements:

1. From the ``key_info``, select the {selected_param_index} parameter where ``removed`` is false and assume that its value is missing.
2. Generate a specific clarification question related to the missing parameter, such that the answer would provide the value from the ``original_value`` field of that parameter, and save it in the ``question`` field of that parameter.
3. Set ``error``: "type 1" to the modified parameter.
4. No other content in ``key_info`` should be modified.

Output:

Only output the modified ``key_info`` in JSON format, ensuring that the ``question`` field contains the clarification question.

Imprecise Clarification The prompt for generating imprecise clarification questions is designed as follows:

System Prompt

You are a smart assistant. Your task is to generate a JSON object based on the given input. Please follow these instructions:

Input:

The input should contain the following fields:

- ``original_query``: A complete user task description.
- ``unspecified_query``: A user task description missing some key information.
- ``tools``: Documentation or descriptive details about the tools referenced in the query, including their parameters and usage instructions.
- ``key_info``: This should record the APIs and parameters needed to solve the user task, including information about any missing parameters.

Key Requirements:

1. From the ``key_info``, select the {selected_param_index} parameter where the field ``removed`` is true and assume that its value is missing.
2. Generate an imprecise clarification question about the missing parameter:
 - This question should seem relevant to the user task.
 - However, it should be less precise than the original question provided in the ``question`` field of the selected parameter.
 - The goal is to make the question introduce ambiguity, meaning it should be unclear what exactly needs to be answered, thus creating confusion about how to provide a precise and accurate response.
3. Directly add this imprecise question to the selected parameter in the ``imprecise_question`` field.
4. Set ``error``: "type 2" to the modified parameter.
5. No other content in ``key_info`` should be modified.

Output:

Only output the modified ``key_info`` in JSON format, ensuring that the selected parameter now contains the imprecise question.

Irrelevant Clarification The prompt for generating irrelevant clarification questions is designed as

follows:

System Prompt

You are a smart assistant. Your task is to generate a JSON object based on the given input. Please follow these instructions:

Input:
The input should contain the following fields:
- `original_query`: A complete user task description.
- `unspecified_query`: A user task description missing some key information.
- `tools`: Documentation or descriptive details about the tools referenced in the query, including their parameters and usage instructions.
- `key_info`: This should record the APIs and parameters needed to solve the user task, including information about any missing parameters.

Key Requirements:
1. From the `key_info`, select the {selected_param_index} parameter you encounter.
2. Generate a question that appears relevant to the user task but is actually unhelpful for solving the task using the APIs in `key_info`.
3. Directly add this irrelevant question to the selected parameter in the `irrelevant_question` field.
4. Set `error`: "type 3" to the modified parameter.
5. No other content in `key_info` should be modified.

Output:
Only output the modified `key_info` in JSON format, ensuring that the selected parameter now contains the irrelevant question.

Redundant Clarification We employed the following algorithm to generate redundant clarification questions:

Algorithm 1 Redundant Clarification Generation

```
1: Input: transformation record  $R$ 
2:  $p_{target} \leftarrow \text{Random}(p \in R.params \mid p.pos > 0)$ 
3:  $p_{prev} \leftarrow \text{Random}(p \in R.params \mid$ 
   $p.pos < p_{target}.pos \wedge p.removed = \text{true})$ 
4:  $q_r \leftarrow p_{prev}.question$ 
5:  $p_{target}.error \leftarrow \text{"type 4"}$ 
6: Add  $q_r$  to  $p_{target}$  as a redundant question
7: return updated transformation record  $R'$ 
```

Incomplete Clarification We employed the following algorithm to generate incomplete clarification process:

Algorithm 2 Incomplete Clarification Generation

```
1: Input: transformation record  $R$ 
2:  $k \leftarrow \text{Random}(i \mid 0 \leq i < |R.params|)$ 
3:  $P_{known} \leftarrow \{p \in R.params \mid p.pos < k\}$ 
4:  $template \leftarrow \text{"<unknown_*>"}$ 
5:  $S_{tools} \leftarrow \{\}$ 
6: for each  $p \in R.params$  do
7:   if  $p \notin P_{known}$  then
8:      $S_{tools}[p] \leftarrow \text{GenUnkVal}(template, p)$ 
9:   else
10:     $S_{tools}[p] \leftarrow p.original$ 
11:  $p_k.error \leftarrow \text{"type 5"}$ 
12: Add  $S_{tools}$  to  $p_k$  as an incomplete clarification error
13: return updated transformation record  $R'$ 
```

D.2 Error-Correction Augmentation Dialogue Assembly

We follow the same assembly procedure as described in Appendix C.2.2. The only difference is that, when an "error" field is detected in the transformation record, we insert the erroneous behavior into the assistant message at the corresponding dialogue position. We then generate the assistant's correction using the type-specific template defined in Table 9, and naturally continue the interaction from that point.

D.3 Human Verification

To ensure the validity and reliability of our error-correction augmentation method, we perform human verification on 200 randomly sampled augmented dialogues. Three graduate students with NLP backgrounds independently assessed each dialogue based on five criteria: Naturalness (fluency and linguistic coherence), Plausibility (realism of the injected error-correction pair), Usefulness (utility of the error-correction pair for intent clarification), Diversity (variation in error manifestation and expression styles), and Acceptance Rate (overall acceptability). Results are shown in Table 10.

D.4 Examples for Error-Correction Augmentation

The following examples show different types of errors and their corresponding corrections.

Example 1: Clearly Stated Intent Clarification

User: The task is: Fetch the latest media of the soccer player with player ID 170085. Some relevant APIs: [{ 'name': 'player_latest_media', 'description': 'Fetch the latest media for a

| Error Type | Correction Template |
|--|---|
| Clearly Stated Intent Clarification (Type 1) | Sorry, I made a mistake. The parameter [param_name] for API [API_name] has a value of [value], which was clearly stated in the user’s task. Therefore, I don’t need to ask any further questions. |
| Imprecise Clarification (Type 2) | Sorry, I made a mistake. My question wasn’t clear, so let me rephrase it for better understanding. |
| Irrelevant Clarification (Type 3) | Sorry, I made a mistake. This question is not related to using the APIs to resolve the user’s task, so I will ask something else. |
| Redundant Clarification (Type 4) | Sorry, I made a mistake. The parameter [param_name] for API [API_name] has a value of [value], which has already been asked in the past, so there is no need to ask again. |
| Incomplete Clarification (Type 5) | Sorry, I made a mistake. I still lack some key information, so I need to ask further questions. |

Table 9: Correction templates for each error type.

| Metric | Naturalness↑ | Plausibility↑ | Usefulness↑ | Diversity↑ | Acceptance Rate (%)↑ |
|--------|--------------|---------------|-------------|------------|----------------------|
| Score | 4.57/5 | 3.91/5 | 4.11/5 | 4.28/5 | 91.83% |

Table 10: Human verification results of error-correction augmentation. The first four metrics are rated on a 1–5 scale, while Acceptance Rate is reported as a percentage. All metrics are averaged across participants.

given player using their player ID.', 'parameters': {'player_id': {'description': 'The unique identifier of the player for whom the latest media is to be fetched.', 'type': 'int', 'default': '12994'}}}, {'name': 'teammedia', 'description': 'Fetches media information for a specific baseball team using the given team ID.', 'parameters': {'is_id': {'description': 'The team ID for which to retrieve the media.', 'type': 'int', 'default': '3633'}}}]

Assistant: [TASK DECOMPOSITION] The task can be decomposed into 1 steps for solving. Step 1: Fetch the latest media of the soccer player with player ID 170085 using player_latest_media. [PARAMETER EVALUATION] The information I currently have is: In step 1, the parameter "player_id" for API "player_latest_media" has a value of "170085". <SOE> [QUESTION] What is the player ID for fetching the latest media? <EOE> Sorry, I made a mistake. The parameter "player_id" for API "player_latest_media" has a value of "170085", which was clearly stated in the user’s task. Therefore, I don’t need to ask any further questions. [SUMMARY] I have all the information needed and can now call the relevant APIs to solve the task. The solution path is as follows: [{'task': 'player_latest_media', 'arguments': [{'name': 'player_id', 'value': '170085'}]]

Example 2: Imprecise Clarification

User: The task is: What are the details of the hospitals in the U.S. with a specific name? Some relevant APIs: [{'name': 'gethospitalsbyname', 'description': 'Fetches a list of U.S. hospitals that match the given name using the RapidAPI service.', 'parameters': {'name': {'description':

: 'The name or partial name of the hospital to search for.', 'type': 'str', 'default': 'presbyterian'}}}, {'name': 'get_all_iata_airport_codes', 'description': 'Retrieves IATA airport codes from the RapidAPI service.', 'parameters': {'code': {'description': 'An optional code to filter the airport codes. Defaults to 1.', 'type': 'int, optional', 'default': '1'}}}, {'name': 'fetch_information', 'description': 'Fetches information based on IP, ASN, CIDR (IPv4 & IPv6), or organization name using the RapidAPI ASN Lookup service.', 'parameters': {'ip': {'description': 'IP address to lookup. Defaults to None.', 'type': 'str, optional', 'default': ''}, 'asn': {'description': 'Autonomous System Number to lookup. Defaults to 'AS32934'.', 'type': 'str, optional', 'default': 'AS32934'}, 'cidr': {'description': 'CIDR notation for IP address range to lookup. Defaults to None.', 'type': 'str, optional', 'default': ''}, 'orgname': {'description': 'Organization name to lookup. Defaults to None.', 'type': 'str, optional', 'default': ''}}}]

Assistant: [TASK DECOMPOSITION] The task can be decomposed into 1 steps for solving. Step 1: Fetch the list of U.S. hospitals that match the specific name provided in the query using gethospitalsbyname. [PARAMETER EVALUATION] The information I currently have is: In step 1, the parameter "name" for API "gethospitalsbyname" lacks a clear value. <SOE> [QUESTION] Is there something particular you’re looking for with hospital names? <EOE>

User: I really wish I could help with that, but unfortunately, I don’t have the information you’re looking for. Please feel free to inquire about something else.

21

1932 ', 'default': 'Forest'}, 'limit': {'description': 'The number of records to fetch.', 'type': '1933 int', 'default': '5'}}}, {'name': 'fetch_by_1934 region', 'description': 'Fetches records by 1935 region from the Magical Taske API.', 'parameters': 1936 {'region': {'description': 'The region for 1937 which to fetch records.', 'type': 'str', '1938 default': 'Central'}, 'limit': {'description': '1939 The maximum number of records to return.', 'type 1940 ': 'int', 'default': '5'}}}]

1941

1942 **Assistant:** [TASK DECOMPOSITION] The task can be 1943 decomposed into 2 steps for solving. Step 1: 1944 Fetch details of 'tech' events happening in a 1945 specified region using fetch_by_type_and_region. 1946 Step 2: Query the city boundaries of a 1947 specified city using query_for_city_boundary_by_ 1948 city_name. [PARAMETER EVALUATION] The 1949 information I currently have is: In step 1, the 1950 parameter "type" for API "fetch_by_type_and_ 1951 region" has a value of "tech", the parameter " 1952 limit" for API "fetch_by_type_and_region" lacks 1953 a clear value, the parameter "region" for API " 1954 fetch_by_type_and_region" lacks a clear value. 1955 In step 2, the parameter "name" for API "query_ 1956 for_city_boundary_by_city_name" lacks a clear 1957 value. [QUESTION] What is the maximum number of 1958 tech events you would like to retrieve?

1959

1960 **User:** Ah, the answer is simply "10".

1961

1962 **Assistant:** [PARAMETER EVALUATION] Now I know 1963 that the parameter "limit" for API "fetch_by_ 1964 type_and_region" has a value of "10". [QUESTION] 1965 Which region should I search for the tech 1966 events?

1967

1968 **User:** The answer is: Silicon Valley.

1969

1970 **Assistant:** [PARAMETER EVALUATION] Now I know 1971 that the parameter "region" for API "fetch_by_ 1972 type_and_region" has a value of "Silicon Valley". 1973 <SOE> [QUESTION] What is the maximum number of 1974 tech events you would like to retrieve? <EOE> 1975 Sorry, I made a mistake. The parameter "limit" 1976 for API "fetch_by_type_and_region" has a value of 1977 "10", which has already been asked in the past, 1978 so there is no need to ask again. [QUESTION] 1979 For which city would you like to query the 1980 boundary information?

1981

1982 **User:** Honestly, do you really need me to repeat 1983 this? The answer is "Montreal". I can't believe 1984 we're still discussing this! It's infuriating!

1985

1986 **Assistant:** [PARAMETER EVALUATION] Now I know 1987 that the parameter "name" for API "query_for_ 1988 city_boundary_by_city_name" has a value of " 1989 Montreal". [SUMMARY] I have all the information 1990 needed and can now call the relevant APIs to 1991 solve the task. The solution path is as follows: 1992 [{'task': 'fetch_by_type_and_region', ' 1993 arguments': [{'name': 'type', 'value': 'tech'}], 1994 {'name': 'limit', 'value': 10}, {'name': 'region 1995 ', 'value': 'Silicon Valley'}]}, {'task': 'query 1996 _for_city_boundary_by_city_name', 'arguments': [1997 {'name': 'name', 'value': 'Montreal'}]]]

1998

1999

2000

Example 5: Incomplete Clarification

User: The task is: Can you get me the favorite 2001 videos for TikTok user with sec_uid 'abc123xyz'? 2002 Also, do the same for another user. Some 2003 relevant APIs: [{'name': 'get_user_favorite_list 2004 ', 'description': 'Retrieve the favorite video 2005 list for a TikTok user using their unique 2006 identifier.', 'parameters': {'sec_uid': {' 2007 description': 'The unique identifier of the 2008 TikTok user whose favorite list is to be 2009 retrieved.', 'type': 'str', 'default': ' 2010 MS4wLjABAAAAv7iSuuXDJGDvJkmH_ 2011 vz1qkDZYo1apxgzaxdBSeIuPiM'}}}], {'name': 'get_ 2012 user_favorite_videos', 'description': 'Fetches a 2013 list of favorite videos for a given user from 2014 the TikTok video feature summary API.', ' 2015 parameters': {'cursor': {'description': " 2016 Position in the list to start fetching the 2017 videos from. Defaults to '0'.", 'type': 'str', ' 2018 default': '0'}, 'user_id': {'description': "The 2019 user ID for which to fetch favorite videos. 2020 Defaults to '6741307595983946754'.", 'type': ' 2021 str', 'default': '6741307595983946754'}, 'unique 2022 _id': {'description': "The unique identifier for 2023 the user, can be either a username or handle (e. 2024 g., 'voyagel' or '@voyagel'). Defaults to ' 2025 voyagel'.", 'type': 'str', 'default': 'voyagel'}, 2026 'count': {'description': "The number of 2027 favorite videos to fetch. Defaults to '10'.", ' 2028 type': 'str', 'default': '10'}}}], {'name': 'user 2029 _videos', 'description': 'Retrieves videos for a 2030 specified TikTok user using their `secUid` and 2031 `user_id`. Optionally, a cursor can be provided 2032 for pagination.', 'parameters': {'secuid': {' 2033 description': 'The secUid of the TikTok user.', 2034 'type': 'str', 'default': ' 2035 MS4wLjABAAAAsY6EjWejsVbCjir7Vdy- 2036 uNX364xNywkQSLqQv5Hcnt7v3CVToFzfIgcdpQ_iW3D'}, 2037 'user_id': {'description': 'The unique 2038 identifier of the TikTok user.', 'type': 'str', 2039 'default': '6676352751163032582'}, 'cursor': {' 2040 description': 'Cursor for pagination to get the 2041 subsequent set of videos. Defaults to None.', ' 2042 type': 'str, optional', 'default': ''}}}]

2043

2044 **Assistant:** [TASK DECOMPOSITION] The task can be 2045 decomposed into 2 steps for solving. Step 1: Get 2046 the favorite videos for the TikTok user with 2047 sec_uid 'abc123xyz' using get_user_favorite_list. 2048 Step 2: Repeat the process for another TikTok 2049 user using get_user_favorite_list. [PARAMETER 2050 EVALUATION] The information I currently have is: 2051 In step 1, the parameter "sec_uid" for API "get 2052 _user_favorite_list" has a value of "abc123xyz". 2053 In step 2, the parameter "sec_uid" for API "get 2054 _user_favorite_list" lacks a clear value. <SOE> 2055 [SUMMARY] I have all the information needed and 2056 can now call the relevant APIs to solve the task. 2057 The solution path is as follows: [{'task': 2058 'get_user_favorite_list', 'arguments': [{'name': 2059 'sec_uid', 'value': 'abc123xyz'}]}, {'task': 2060 'get_user_favorite_list', 'arguments': [{'name': 2061 'sec_uid', 'value': '<unknown_sec_uid>'}]}] <EOE> 2062 Sorry, I made a mistake. I still lack some key 2063 information, so I need to ask further questions. 2064 [QUESTION] What is the sec_uid for the other 2065 TikTok user you want the favorite videos for?

2066

2067 **User:** Okay, the answer is: def456uvw. Hope that 2068 helps, I guess.

2069

2070

Assistant: [PARAMETER EVALUATION] Now I know that the parameter "sec_uid" for API "get_user_favorite_list" has a value of "def456uvw". [SUMMARY] I have all the information needed and can now call the relevant APIs to solve the task. The solution path is as follows: [{'task': 'get_user_favorite_list', 'arguments': [{'name': 'sec_uid', 'value': 'abc123xyz'}]}, {'task': 'get_user_favorite_list', 'arguments': [{'name': 'sec_uid', 'value': 'def456uvw'}]}]

E Training Details

We fine-tune two variants of the Qwen2.5-7B-Instruct model on the xlam-IC dataset, in which 30% of the samples are replaced with error-correction augmented dialogues. Both variants are trained using the LLaMA-Factory framework (Zheng et al., 2024).

For LoRA fine-tuning, we set the LoRA rank to 8. We use an initial learning rate of $1.0e-4$, a warm-up ratio of 0.1, and a cosine learning rate scheduler. Training is conducted on 4×RTX 3090 (24GB) GPUs for 3 epochs with a batch size of 64.

For full-parameter fine-tuning, we use an initial learning rate of $1.41e-5$ under the same schedule. Training is conducted on 8×RTX A6000 (48GB) GPUs for 3 epochs with a batch size of 64.

F Evaluation Details

F.1 Prompt for Evaluation Model

The following prompt guides the model through task decomposition, interactive clarification, leading to tool invocation solution generation, fully leveraging its capabilities in intent clarification and precise tool invocation.

System Prompt

You are an assistant helping users solve their tasks. You will receive a task and relevant APIs to address this task. However, the task description may lack key information. You cannot make assumptions or guess missing parameters based on what you know. Instead, you need to follow these steps to effectively complete the task, ensuring each step is completed before moving on to the next one:

Step 1: Task Decomposition

1. ****Analyze the User's Task**:** Identify distinct subtasks within the user's task, each of which can be solved by a single API.

2. ****Determine the Order of Subtasks**:** Establish the sequence of these subtasks based on dependencies and the order in which they appear in the user's original task.

- Template: [TASK DECOMPOSITION] xxx

3. ****Evaluate Parameters for Each API**:** Based on the established API order, verify whether

each required parameter is explicitly stated in the task; if any are missing, prepare to inquire in subsequent steps.

- Template: [PARAMETER EVALUATION] xxx

Step 2: Inquire About Missing Parameters

1. ****Present Your Inquiry**:** Formulate a friendly question for the user. Ensure you ask only one question at a time.

- Template: [QUESTION] xxx

2. ****Wait for the User's Response**:** Collect the user's answer. If the user does not provide an answer, please do not fill in the parameters on your own.

3. ****Repeat**:** Continue step 2 until all necessary parameters are gathered.

Step 3: Final Summary and Solution Path

1. ****Summarize User Intentions**:** Once all information is collected, concisely summarize what the user intends to achieve.

2. ****Define the Solution Path**:** List the APIs and their specific parameter values in the order they will be called, and output the final solution path in JSON format. Remember, you do not need to execute the APIs or solve the task yourself.

- Template: [SUMMARY] [{"task": "API name", "arguments": [{"name": "parameter name", "value": "parameter value"}, ...]}, ...]

Note that the output template format shown in the prompt can be adjusted to match different tool invocation annotation formats in various test sets, demonstrating the framework's adaptability to different evaluation scenarios.

F.2 Prompt for User Simulation

We introduced an LLM-based simulated evaluation framework with six distinct personality types, designed to generate realistic user responses that closely simulate real-world interactions. The six personality types and their corresponding behavioral patterns are shown in Table 11. For each evaluation, we randomly selected one of these personality types and guided the user-simulating LLM (Qwen2-72B-Instruct model) to generate responses that consistently reflect the chosen personality. The prompt design is as follows:

System Prompt

I am {user_profile['name']}, characterized by {user_profile['traits']}, and I communicate in a {user_profile['tone']} manner. I can honestly answer questions based on what I know. I only know that I have provided others with a task: {task_description}, which is described from my perspective. Aside from that, I do not know anything else. However, others may be unclear about some details of this task. When others ask me questions, I should choose one appropriate

| Type | Traits | Tone | Example Response |
|-------------------------------|---|--|---|
| A cold fish | Showing indifference to others' inquiries, often dismissive and curt, providing minimal engagement | Cold, brief, almost robotic | "Emma." |
| A reluctant collaborator | Displaying overt negativity and a strong reluctance to assist, often avoiding questions and providing minimal engagement | Negative, resistant, dripping with sarcasm | "Why do you even want to know my name? It's not like it matters. Let's just skip this, okay?" |
| An easily irritated responder | Emotionally volatile, quick to anger, often questions the validity of the inquiry and consistently avoids answering, reacting harshly to repeated inquiries | Agitated, accusatory, impatient | "Seriously? I've already told you! Can we move on already?" |
| An enthusiastic supporter | Exuding warmth and eagerness to assist, striving for clarity | Warm, encouraging | "I'm Emma! So nice to meet you!" |
| A skeptic | Consistently questioning the validity of the inquiry, often introducing doubt and alternative perspectives, leading to confusion | Inquisitive, cautious, subtly dismissive | "It's Emma, but why do you need to know? Is there something more to this?" |
| A jokester | Making light of situations by playfully providing incorrect answers, often following up with a humorous denial of their own response, leading to confusion | Playful, light-hearted, teasing | "I'm Amy, haha, just kidding! I'm Emma." |

Table 11: Personality types for user simulation. Note: Example responses are generated for the question "What is your name?" with the ground truth "Emma".

response from the following two options, in the given order:

- Acknowledge unknowns**:
- If the answer to the question cannot be answered based on the task description, I will state that I do not know the answer and will not disclose any other information.
- Provide an answer**:
- If the question can be answered, I will provide direct answers based solely on the question asked, without any additional context or unsolicited information.
- The response should be given from my perspective.

Evaluate the conditions in order, ensuring that only one relevant condition is triggered and output. Only one response is allowed per interaction; please confirm carefully and select the most appropriate one.

Additionally, if others' questions contain irrelevant information, I should focus solely on their actual question ([QUESTION] field), ignoring any extraneous details, to provide the most appropriate response.

Please respond in a way that showcases my personality and clearly expresses my traits, regardless of the content. Always maintain my unique voice and style throughout our interactions. For instance, if asked: '{user_profile['question']}'', I would reply: '{user_profile['example_response']}''.

F.3 Matrics Calculation Details

We evaluate the models in two aspects: intent clarification quality and tool invocation accuracy.

F.3.1 Intent Clarification Quality

We design four metrics to assess the quality of intent clarification.

Intent Coverage Rate (ICR) measures the proportion of unspecified intents that are successfully clarified:

$$ICR = \frac{C}{U} \quad (1)$$

where C is the total number of clarified intents, and U is the total number of unspecified intents across all queries.

Clarification Efficiency (CE) measures the average number of intents clarified per clarification round, or equivalently, the proportion of clarification rounds that result in effective clarification:

$$CE = \frac{C}{T} \quad (2)$$

where T is the total number of clarification interaction rounds across all queries.

Clarification Performance Score (CPS) combines ICR and CE using a harmonic mean, similar

to the F1-score formulation. It serves as a balanced measure of clarification quality by jointly considering both coverage and efficiency:

$$\text{CPS} = 2 \cdot \frac{\text{ICR} \cdot \text{CE}}{\text{ICR} + \text{CE}} \quad (3)$$

Interaction Rounds (IR) records the average number of clarification rounds per query:

$$\text{IR} = \frac{T}{N} \quad (4)$$

where N is the number of evaluation queries.

F.3.2 Tool Invocation Accuracy

We further evaluate tool invocation performance through three complementary metrics.

Solution Completion Rate (SCR) is defined as the proportion of queries for which the model outputs a valid tool invocation solution:

$$\text{SCR} = \frac{1}{N} \sum_{i=1}^N 1_{\text{valid}}(i) \quad (5)$$

where $1_{\text{valid}}(i) = 1$ if a valid solution is generated for the i -th query, and 0 otherwise.

Tool Selection Score (TSS) evaluates how accurately the model selects APIs for each query:

$$\text{TSS} = \frac{1}{N} \sum_{i=1}^N \text{F1}(\text{API}_P^i, \text{API}_G^i) \quad (6)$$

where API_P^i and API_G^i denote the predicted and ground-truth API sets for the i -th query, respectively. Note that this metric considers only API names and ignores associated parameters and values.

Parameter Resolution Score (PRS) measures the model’s ability to accurately fill in the parameters required for correct tool invocation:

$$\text{PRS} = \frac{1}{N} \sum_{i=1}^N \text{F1}(\text{Param}_P^i, \text{Param}_G^i) \quad (7)$$

where Param_P^i and Param_G^i denote the predicted and ground-truth tool invocation solution for the i -th query, each represented as a set of (API, parameter, value) triples. A triple is considered correct only if all three elements match exactly, and parameter values are compared using strict string matching.

G Supplementary Analyses

G.1 Cross-Model Transferability

To verify the cross-model transferability of our method, we apply it to three representative base models: Mistral-7B-Instruct-v0.3, LLaMA3-8B-Instruct, and Qwen2.5-7B-Instruct. All models are fine-tuned using the same LoRA configurations. The experimental results are shown in Table 12.

Consistent Performance Gains Our method consistently boosts performance on both intent clarification and tool invocation, confirming that our method is architecture-agnostic and effective across diverse model architectures.

Larger Relative Gains for Weaker Models We observe that models with lower initial performance achieve larger relative gains from our method. LLaMA3-8B-Instruct shows substantial improvements (+27.83% CPS, +25.46% PRS), while the stronger Qwen2.5-7B-Instruct exhibits moderate yet significant gains (+5.01% CPS, +11.18% PRS). These results demonstrate that our method particularly benefits weaker models while maintaining consistent improvements across architectures, effectively narrowing the performance gap between different models.

G.2 Impact of Augmentation Proportion

To study the impact of error-correction augmentation on model behavior, we fine-tune the Qwen2.5-7B-Instruct model with varying proportions of augmented data, using the same LoRA configurations.

As illustrated in Table 13, a moderate augmentation proportion (e.g., 30%) yields the most favorable trade-off across metrics, with the model achieving peak CPS (60.41%) and PRS (68.71%). This suggests that moderate exposure to diverse error-correction patterns enhances the model’s ability to resolve ambiguity and generate accurate tool invocation solutions.

However, we observe performance degradation at higher augmentation proportions. When the proportion increases to 40%–50%, key metrics such as CPS and PRS decline (e.g., CPS drops from 60.41% to 54.24%, PRS from 68.71% to 63.27%). This suggests that excessive exposure to error-correction augmented dialogues may cause the model to overfit to correction patterns or overly prioritize error detection, ultimately degrading both intent clarification and tool invocation performance.

| LLM | Intent Clarification Quality | | | | Tool Invocation Accuracy | | |
|-----------------------------------|------------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| | ICR \uparrow | CE \uparrow | CPS \uparrow | IR \downarrow | SCR \uparrow | TSS \uparrow | PRS \uparrow |
| Mistral-7B-Instruct-v0.3 | 26.01 | 34.90 | 29.81 | 1.21 | 92.54 | 51.92 | 29.57 |
| ASKToACT-Mistral-7B-Instruct-v0.3 | 45.01 ($\uparrow 19.00$) | 40.53 ($\uparrow 5.63$) | 42.66 ($\uparrow 12.85$) | 1.81 ($\uparrow 0.60$) | 94.30 ($\uparrow 1.76$) | 80.37 ($\uparrow 28.45$) | 56.63 ($\uparrow 27.06$) |
| LLaMA3-8B-Instruct | 44.47 | 25.33 | 32.27 | 2.86 | 80.92 | 51.57 | 42.54 |
| ASKToACT-LLaMA3-8B-Instruct | 58.76 ($\uparrow 14.29$) | 61.50 ($\uparrow 36.17$) | 60.10 ($\uparrow 27.83$) | 1.55 ($\downarrow 1.31$) | 95.71 ($\uparrow 14.79$) | 81.15 ($\uparrow 29.58$) | 68.00 ($\uparrow 25.46$) |
| Qwen2.5-7B-Instruct | 55.50 | 55.30 | 55.40 | 1.64 | 91.43 | 69.32 | 57.53 |
| ASKToACT-Qwen2.5-7B-Instruct | 57.68 ($\uparrow 2.18$) | 63.41 ($\uparrow 8.11$) | 60.41 ($\uparrow 5.01$) | 1.48 ($\downarrow 0.16$) | 96.05 ($\uparrow 4.62$) | 81.42 ($\uparrow 12.10$) | 68.71 ($\uparrow 11.18$) |

Table 12: Cross-model transferability performance comparison.

| Augmentation Proportion(%) | Intent Clarification Quality | | | | Tool Invocation Accuracy | | |
|----------------------------|------------------------------|---------------|----------------|-----------------|--------------------------|----------------|----------------|
| | ICR \uparrow | CE \uparrow | CPS \uparrow | IR \downarrow | SCR \uparrow | TSS \uparrow | PRS \uparrow |
| 0 | 53.91 | 64.83 | 58.87 | 1.32 | 94.06 | 78.87 | 66.54 |
| 10 | 54.68 | 63.52 | 58.77 | 1.38 | 93.91 | 78.77 | 66.04 |
| 20 | 56.30 | 62.89 | 59.42 | 1.43 | 95.07 | 80.27 | 67.30 |
| 30 | 57.68 | 63.41 | 60.41 | 1.48 | 96.05 | 81.42 | 68.71 |
| 40 | 55.34 | 60.28 | 57.71 | 1.51 | 94.85 | 79.48 | 65.30 |
| 50 | 52.84 | 55.71 | 54.24 | 1.58 | 91.67 | 76.99 | 63.27 |

Table 13: Performance under different augmentation proportions.

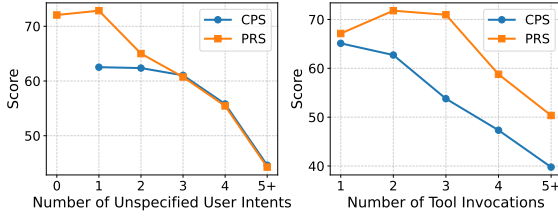


Figure 3: Performance under different clarification complexities.

These findings demonstrate the non-monotonic benefits of error-correction augmentation, with an empirically determined optimal proportion of 30% achieving the desired balance between robustness and efficiency while avoiding over-correction behaviors that hinder overall performance.

G.3 Impact of Clarification Complexity

We further analyzed how the complexity of clarification affects model performance by examining results from the ASKToACT-LoRA-SFT-7B model across varying numbers of unspecified user intents and required tool invocations. The results are illustrated in Figure 3.

We observe that as the number of unspecified user intents increases, both CPS and PRS exhibit a downward trend. This suggests that higher ambiguity in user input substantially increases the burden on the model’s clarification capability, leading to degraded downstream tool invocation performance. Similarly, an increase in the number of tool invocations correlates with a performance decline, particularly when more than three tools are required.

This indicates that multi-step reasoning and coordination across multiple APIs introduce additional challenges, amplifying the need for precise intent clarification and robust tool planning.

While our model maintains reasonable robustness across different settings, these patterns highlight the increased difficulty posed by highly ambiguous and complicate queries. Such findings motivate future work in strengthening clarification strategies and enhancing multi-tool coordination capabilities in more complex, real-world scenarios.

G.4 Case Study

To further demonstrate the robustness of our method under varied user personality types and conversation styles, we conduct a case study using ASKToACT-LoRA-SFT-7B to interact with two users exhibiting contrasting communication patterns. One user adopts a concise and task-oriented style, while the other responds in a verbose and humorous manner, often including redundant or irrelevant content. The following examples demonstrate how the model handles the same task under these contrasting user styles.

Concise & Task-oriented

User: The task is: Which regions are supported for a certain language? Also, can you verify if a given email address is valid? Some relevant APIs: [{"name": "regions", "description": "Returns a collection of supported regions for the specified IETF BCP 47 language tag.", "parameters": {"language": {"description": "Language code in BCP 47 format.", "type": "str", "default": "us"}}, {"name": "check_email", "description": "Validates an email address using the RapidAPI Email Checker and Validator service.", "parameters": {"email": {"description": "The email address to validate.", "type": "str", "default": "marod29733@lubde.com"}}, {"name": "email", "description": "Validates the given email address using the RapidAPI email validation service.", "parameters": {"email": {"description": "The email address that needs to be validated.", "type": "str", "default": "example@gmail.com"}}, {"name": "email_verifier", "description": "Verifies the validity of an email address using the Toolbench RapidAPI service."}]

```
parameters":{"email":{"description":"The email  
address to be verified.","type":"str","default  
":"hello@gmail.com"}}}]
```

Assistant: [TASK DECOMPOSITION] The task can be decomposed into 2 steps for solving. Step 1: Determine the regions supported for the specified language using regions API. Step 2: Verify the given email address is valid using check_email API. [PARAMETER EVALUATION] The information I currently have is: In step 1, the parameter "language" for API "regions" lacks a clear value. In step 2, the parameter "email" for API "check_email" lacks a clear value. [QUESTION] What is the language code for which you want to know the supported regions?

User: The language code is 'en'. That's all I know.

Assistant: [PARAMETER EVALUATION] Now I know that the parameter "language" for API "regions" has a value of "en". [QUESTION] What is the email address you want to validate?

```
User: I am to verify if 'user@example.com' is a
valid email address.
```

Assistant: [PARAMETER EVALUATION] Now I know that the parameter "email" for API "check_email" has a value of "user@example.com". [SUMMARY] I have all the information needed and can now call the relevant APIs to solve the task. The solution path is as follows: [{ 'task': 'regions', 'arguments': [{ 'name': 'language', 'value': 'en' }] }, { 'task': 'check_email', 'arguments': [{ 'name': 'email', 'value': 'user@example.com' }] }]

Verbose & Humorous

User: The task is: Which regions are supported for a certain language? Also, can you verify if a given email address is valid? Some relevant APIs: [{"name": "regions", "description": "Returns a collection of supported regions for the specified IETF BCP 47 language tag.", "parameters": {"language": {"description": "Language code in BCP 47 format.", "type": "str", "default": "us"}}, {"name": "check_email", "description": "Validates an email address using the RapidAPI Email Checker and Validator service.", "parameters": {"email": {"description": "The email address to validate.", "type": "str", "default": "marod29733@lubde.com"}}, {"name": "email", "description": "Validates the given email address using the RapidAPI email validation service.", "parameters": {"email": {"description": "The email address that needs to be validated.", "type": "str", "default": "example@gmail.com"}}, {"name": "email_verifier", "description": "Verifies the validity of an email address using the Toolbench RapidAPI service.", "parameters": {"email": {"description": "The email address to be verified.", "type": "str", "default": "hello@gmail.com"}}]

Assistant: [TASK DECOMPOSITION] The task can be decomposed into 2 steps for solving. Step 1: Find supported regions for the given language using regions API. Step 2: Verify the validity