

# STRUCTURED UNCERTAINTY GUIDED CLARIFICATION FOR LLM AGENTS

Manan Suri<sup>♣</sup>, Puneet Mathur<sup>◇</sup>, Nedim Lipka<sup>◇</sup>

Franck Dernoncourt<sup>◇</sup>, Ryan A. Rossi<sup>◇</sup>, Dinesh Manocha<sup>♣</sup>

<sup>♣</sup>University of Maryland, College Park <sup>◇</sup>Adobe Research

manans@umd.edu, puneetm@adobe.com

## ABSTRACT

LLM agents extend large language models with tool-calling capabilities, but ambiguous user instructions often lead to incorrect invocations and task failures. We introduce a principled formulation of structured uncertainty over tool-call parameters, modeling joint tool-argument clarification as a POMDP with Expected Value of Perfect Information (EVPI) objective for optimal question selection and aspect-based cost modeling to prevent redundancy. Our SAGE-Agent leverages this structured uncertainty to achieve superior efficiency: increasing coverage on ambiguous tasks by 7–39% while reducing clarification questions by 1.5–2.7 $\times$  compared to strong prompting and uncertainty-based baselines. We present Clarify-Bench, the first multi-turn tool-augmented disambiguation benchmark with realistic LLM-based user simulation across diverse domains including document editing, vehicle control, and travel booking. Additionally, we demonstrate that structured uncertainty provides effective training signals for reinforcement learning, boosting When2Call accuracy from 36.5% to 65.2% (3B model) and 36.7% to 62.9% (7B model) through uncertainty-weighted GRPO training. These results establish structured uncertainty as a principled, efficient approach for tool-augmented agents, improving both task success and interaction efficiency in real-world scenarios.

## 1 INTRODUCTION

LLM Agents are AI systems that extend large language models (LLMs) with the ability to take real-world actions autonomously accumulate observations Huang et al. (2024b). These agents often invoke external APIs and tools based on structured function definitions, enabling interaction with databases, web services, and software applications Schick et al. (2023). These agents have been successfully deployed across diverse domains including travel planning, document processing, finance, vehicle control, and drug discovery Xie et al. (2024); Mathur et al. (2024); Yu et al. (2024); Huang et al. (2024a); Liu et al. (2024). However, their effectiveness is fundamentally limited by ambiguous or incomplete user instructions that lead to incorrect tool invocations, failed transactions, and degraded user experience—problems that become increasingly critical as these systems handle more complex, high-stakes tasks.

Ambiguity in user requests poses unique challenges for LLM agents, where imprecise interpretation can cascade into costly execution errors Wang et al. (2024); Vijayvargiya et al. (2025). User ambiguity manifests through vague task specifications ("*find me a good restaurant*"), incomplete parameters ("*book a meeting for tomorrow*"), or implicit assumptions about system capabilities Wang et al.

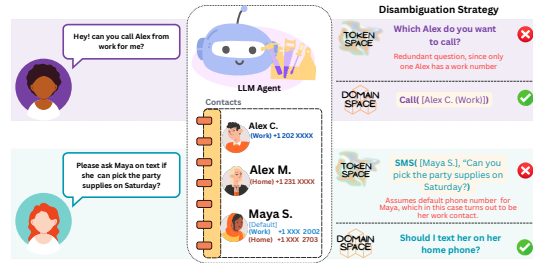


Figure 1: Disambiguation strategies purely grounded in linguistics fail to effectively leverage domain schemas, leading to issues like unnecessary clarifications and assumption of inappropriate default arguments. In contrast, grounding the disambiguation in the structured space of parameter domains mitigates these problems.

(2025). The structured nature of API schemas—with their specific parameter types, constraints, and interdependencies—amplifies this challenge, as a single ambiguous user query often maps to multiple valid API configurations with vastly different outcomes Bandlamudi et al. (2025). For example, *"cancel my subscription"* could apply to multiple services, cancellation types (pause vs. permanent), or effective dates, each requiring different API calls with distinct consequences.

Existing disambiguation approaches suffer from fundamental limitations in the agentic tool-calling context. Due to their next-token prediction training, LLMs often hallucinate missing arguments when faced with incomplete information, leading to incorrect tool invocations Wang et al. (2024). Current methods operate primarily in unstructured language spaces—generating clarifying questions as arbitrary text sequences through prompting strategies—rather than leveraging the structured constraints and dependencies that define tool schemas Kobalczyk et al. (2025); Zhang et al. (2024). While prompting improvements can enhance question phrasing, they cannot fundamentally address the core limitation: without explicit modeling of parameter relationships, importance hierarchies, and feasibility constraints, agents lack principled criteria for determining which questions to ask and when to stop asking them. This results in over-clarification of low-impact details, under-clarification of critical missing information, and inability to distinguish feasible from infeasible requests, as demonstrated in fig. 1.

**Contributions:** ➤ We introduce a principled formulation of *structured uncertainty* over tool-call parameters and model joint tool–argument clarification as a POMDP, using a Bayesian Value of Information objective to optimally select clarification questions that maximize the expected value of perfect information. ➤ SAGE-Agent; Extensive experiments demonstrate that our structured uncertainty–guided strategy substantially improves task success rates while asking the fewest clarification questions, outperforming strong prompting- and uncertainty-based baselines in agentic settings. ➤ We present *ClarifyBench*, the first benchmark dedicated to dynamic, multi-turn tool-calling disambiguation, equipped with an LLM-based user simulator that supports realistic conversational progression and task continuation across diverse domains including document editing, vehicle control, stock trading, travel booking, and file system manipulation. ➤ Finally, we show that our uncertainty formulation serves as an effective reward model, enabling more sample-efficient reinforcement learning fine-tuning for LLM agents in tool-augmented settings.

## 2 RELATED WORK

The challenge of resolving ambiguity in user interaction with LLMs through clarifying questions has gained increasing attention, particularly in tool-calling contexts. Early approaches to clarification focused on general dialogue systems, developing ranking-based methods for question selection Rao & Daumé III (2018); Xu et al. (2019) and Seq2Seq generation Deng et al. (2022). Recent work has specifically addressed ambiguity in tool-calling scenarios: Ask-before-Plan introduces proactive planning agents that predict clarification needs and collect information before execution Zhang et al. (2024), while Active Task Disambiguation frames the problem through Bayesian Experimental Design to maximize information gain from clarifying questions Kobalczyk et al. (2025). Zhang and Choi propose intent-similarity based uncertainty estimation to determine when clarification is beneficial across various NLP tasks Zhang & Choi (2023). Related efforts explore implicit intention understanding in language agents Qian et al. (2024) and proactive dialogue systems that can handle ambiguous queries through goal planning Deng et al. (2023). However, these approaches primarily operate in the general language space without leveraging the structured nature of tool schemas.

## 3 THEORY

Modern LLM agents extend beyond text generation to become *agentic systems* that can interact with external tools and APIs to accomplish complex tasks. These agents typically follow a perception-reasoning-action cycle: they receive user queries, reason about appropriate actions, select and parameterize tool calls, and execute them to achieve desired outcomes. However, this paradigm faces a fundamental challenge when user queries are ambiguous or underspecified—the agent must somehow resolve uncertainty about both *which* tool to use and *how* to parameterize it.

### 3.1 TOOL-CALLING AGENT FRAMEWORK

We model an LLM agent as a system  $\mathcal{M}$  with access to a toolkit  $\mathcal{T} = \{T_1, T_2, \dots, T_K\}$ . Each tool  $T_i$  is characterized by a structured interface that defines its capabilities and parameter requirements.

**Definition 1 (Tool Schema).** A tool  $T_i$  is defined by the tuple  $(name_i, \Theta_i, \mathcal{D}_i, \mathcal{R}_i)$  where:

- $name_i \in \mathbb{S}$  is the tool identifier
- $\Theta_i = \{\theta_{i,1}, \theta_{i,2}, \dots, \theta_{i,m_i}\}$  is the parameter set
- $\mathcal{D}_i = \{\mathcal{D}_{i,1}, \mathcal{D}_{i,2}, \dots, \mathcal{D}_{i,m_i}\}$  where  $\mathcal{D}_{i,j}$  is the domain of parameter  $\theta_{i,j}$
- $\mathcal{R}_i \subseteq \Theta_i$  specifies required parameters

**Definition 2 (Tool Call Candidate).** A tool call candidate  $c_i$  for tool  $T_i$  is a partial function  $c_i : \Theta_i \rightarrow \mathcal{D}_i \cup \{\perp\}$  where  $c_i(\theta_{i,j}) = \perp$  indicates an unspecified parameter.

The agent’s task is to map from an ambiguous natural language query  $q$  to a fully specified tool call  $c^* = (T^*, \theta^*)$  where all required parameters are specified. The *candidate space*  $\mathcal{C} = \{(T_i, c_i) : T_i \in \mathcal{T}, c_i \text{ is valid for } T_i\}$  represents all possible completions consistent with current information.

🔍 **Uncertainty Quantification:** Methods that model uncertainty or disambiguation needs based on LLM response distributions must compute  $p(\text{ambiguous}|u) = \sum_{\mathbf{w}} f(\mathbf{w}) p_{LLM}(\mathbf{w}|u)$  where  $f$  determines if response  $\mathbf{w}$  indicates ambiguity. This conflates model uncertainty with specification uncertainty since the determination function  $f$  itself depends on model capabilities. Our structured approach directly parameterizes  $p(T_i, \theta_i|u)$ , cleanly separating these uncertainty sources.

### 3.2 STRUCTURED BELIEF STATE AND CANDIDATE SCORING

We maintain beliefs over the structured candidate space rather than unstructured text sequences, enabling precise uncertainty quantification.

**Definition 3 (Belief State).** At time  $t$ , we maintain

$$\mathcal{B}(t) = \{(c_i, \Pi_i(t)) : c_i \in \mathcal{C}\}$$

where  $\Pi_i(t) \in [0, 1]$  represents candidate viability. We decompose the joint probability as

$$p(T_i, \theta_i | u) = p(\theta_i | T_i, u) p(T_i | u)$$

and assume a uniform prior over tools, where  $u$  represents the user-input.<sup>1</sup>

Candidate viability then becomes

$$\Pi_i(t) \propto p(\theta_i | T_i, \text{observations}_t) = \prod_{j=1}^{m_i} p(\theta_{i,j} | T_i, \text{observations}_t)$$

where we make a naive independence assumption across parameters for tractability, and parameter certainty is  $p(\theta_{i,j} | T_i, \text{obs}_t) = 1$  if specified,  $|\mathcal{D}_{i,j}(t)|^{-1}$  if unspecified with finite domain, and  $\epsilon$  ( $0 < \epsilon \ll 1$  to approximate uniform density) if the domain is infinite/continuous. Here,  $\mathcal{D}_{i,j}(t)$  is the feasible parameter domain after incorporating observational constraints.

### 3.3 POMDP FORMULATION

We formalize the clarification process as a Partially Observable Markov Decision Process (POMDP), providing theoretical grounding for optimal decision-making under uncertainty.

**State Space:**  $\mathcal{S} = \{(T_i, \theta_i) : T_i \in \mathcal{T}, \theta_i \in \mathcal{D}_i\}$  represents true user intent

**Action Space:**  $\mathcal{A} = \mathcal{Q} \cup \{\text{execute}\}$  where  $\mathcal{Q}$  is the space of clarifying questions

<sup>1</sup>This assumption reflects that, in practice, tools are proposed without strong prior bias (e.g., when suggested by an upstream LLM). Future work could incorporate learned tool usage patterns or contextual priors, but uniform assumptions provide a principled, unbiased baseline and make our formulation cleaner.

**Observation Space:**  $\Omega$  consists of natural language responses to clarifying questions

**Belief Update:** After question  $q$  with response  $r$ , beliefs update through domain constraint propagation:

$$\mathcal{D}_{i,j}(t+1) = \mathcal{D}_{i,j}(t) \cap \text{ExtractConstraints}(r, \theta_{i,j}, T_i) \quad (1)$$

$$\pi_i(t+1) \propto \pi_i(t) \cdot P(r|c_i, q) \cdot \prod_j p(\theta_{i,j}|T_i, \text{observations}_{t+1}) \quad (2)$$

**Reward Function:**  $R = \mathbb{1}[\text{execution matches true intent}] - \lambda \sum_q \text{Cost}(q)$

### 3.4 BAYESIAN VALUE OF INFORMATION FOR QUESTION SELECTION

Drawing from Bayesian Decision Theory and Value of Information frameworks (Rainforth et al., 2024), we formalize optimal question selection through Expected Value of Perfect Information (EVPI).

Combining current expected value  $V_{\text{current}} = \max_{c_i \in \mathcal{C}} \pi_i(t)$  with expected value after questioning  $V_{\text{after}}(q) = \mathbb{E}_{r \sim P(r|q, \mathcal{B}(t))} [\max_{c_i \in \mathcal{C}} \pi_i(t|q, r)]$ , we obtain:

**Definition 4 (Expected Value of Perfect Information).**

$$\text{EVPI}(q, \mathcal{B}(t)) = \mathbb{E}_{r \sim P(r|q, \mathcal{B}(t))} \left[ \max_{c_i \in \mathcal{C}} \pi_i(t|q, r) \right] - \max_{c_i \in \mathcal{C}} \pi_i(t) \quad (3)$$

where response is assumed as  $P(r|q, \mathcal{B}(t)) = \sum_i \pi_i(t) P(r|c_i, q)$ . EVPI naturally handles both tool disambiguation and parameter clarification in a unified framework—questions helping resolve tool choice and parameter values are evaluated using the same information-theoretic criterion.

### 3.5 COST-AWARE SELECTION WITH REDUNDANCY PREVENTION

**Aspects and normalized beliefs.** We introduce the atomic unit of disambiguation, which we call an *aspect*. An **aspect** is a parameter-level identifier  $a_{i,j}$  that refers to parameter  $\theta_{i,j}$  of tool  $T_i$ . The full set of aspects is

$$\mathcal{A} \triangleq \{a_{i,j} \mid i \in [1..K], j \in [1..m_i]\}.$$

A clarifying question targets a subset of aspects: for question  $q$  we write  $\mathcal{A}(q) \subseteq \mathcal{A}$ . For bookkeeping we count how often an aspect has been targeted up to time  $t$  as

$$n_a(t) \triangleq |\{i \leq t : a \in \mathcal{A}(q_i)\}|.$$

💡 **Structured Response Handling:** Past disambiguation methods sample from solution distributions  $p(\text{solution}|q)$  to estimate ambiguity, requiring expensive enumeration over response spaces. Our approach treats responses as domain constraints:  $r \rightsquigarrow \mathcal{D}_{i,j}(t+1) = \mathcal{D}_{i,j}(t) \cap C(r)$  where  $C(r)$  extracts constraints. This reduces integration to finite constraint patterns, enabling exact EVPI computation.

Pure information maximization can lead to excessive questioning. We introduce a cost model that penalizes redundant questions about previously addressed aspects.

**Definition 5 (Redundancy Cost).** For question  $q$  targeting aspects  $\mathcal{A}(q)$ , with aspect history  $n_a(t) = |\{i \leq t : a \in \mathcal{A}(q_i)\}|$ :

$$\text{Cost}(q, t) = \lambda \sum_{a \in \mathcal{A}(q)} n_a(t) \quad (4)$$

**Optimal Question Selection and Termination:**

$$q^*(t) = \arg \max_{q \in \mathcal{Q}} [\text{EVPI}(q, \mathcal{B}(t)) - \text{Cost}(q, t)] \quad (5)$$

$$\text{Stop when: } \max_q [\text{EVPI}(q) - \text{Cost}(q)] < \alpha \cdot \max_i \pi_i(t) \quad (6)$$

The linear cost model captures essential redundancy penalties while remaining computationally tractable, with hyperparameter  $\lambda$  providing control over questioning behavior.

## 4 CLARIFYBENCH

The evaluation of clarification strategies in tool-calling agents requires benchmarks that capture the complexity of real-world user interactions, particularly when dealing with ambiguous or infeasible requests. To address this need, we introduce **ClarifyBench**, a comprehensive benchmark designed to evaluate clarification strategies across diverse domains and query types. As shown in Table 1, existing benchmarks exhibit critical limitations: many lack support for ambiguous and infeasible queries entirely, while those that include such scenarios are limited in scope or domain coverage. Moreover, most benchmarks rely on static evaluation and lack dynamic user simulation capabilities essential for evaluating interactive clarification strategies.

**ClarifyBench** addresses these limitations through dynamic user simulation enabling realistic multi-turn interactions, comprehensive query types (normal, ambiguous, and infeasible), and multi-domain coverage across five distinct domains. Figure 2 illustrates the benchmark design: a user simulator conducts multi-turn interactions with tool-equipped LLM agents, simulating genuine conversational progression where users naturally follow up with related requests after clarification exchanges. Evaluation compares ground truth tool calls with agent-generated actions, providing robust assessment of clarification effectiveness across realistic scenarios.

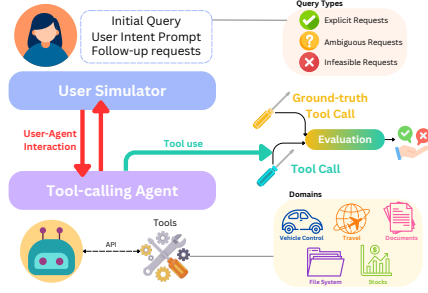


Figure 2: ClarifyBench enables comprehensive evaluation of agent clarification strategies by simulating normal, ambiguous, and infeasible user queries across five domains. A dynamic user simulator conducts multi-turn interactions with tool-equipped LLM agents, with evaluation based on alignment with ground truth agent tool calls.

Benchmark	Dynamic User Simulation	Ambiguous Queries	Infeasible Queries	Multi-turn Requests	Tool Domains	Number of Tools
AgentBoard Ma et al. (2024)	✗	✗	✗	✗	Information Retrieval, Manipulation	50
$\tau$ -bench	✓	✗	✗	✗	Retail, Airlines	24
MMAU Yin et al. (2024)	✗	✗	✗	✗	RapidAPI Tools	364
ToolSandbox Lu et al. (2024)	✓	✗	✗	✗	Personal Assistant	34
Ask-Before-Plan Zhang et al. (2024)	✓	✓	✓	✗	Travel	6
BFCL-v3 Patil et al. (2025)	✗	✓	✗	✓	Vehicle Control, Stocks, Travel, File System	129
<b>ClarifyBench</b>	✓	✓	✓	✓	Documents, Vehicle Control, Stocks, Travel, File System	92

Table 1: Comparison of ClarifyBench with existing tool-calling benchmarks.

### 4.1 BENCHMARK DESIGN

ClarifyBench encompasses five diverse domains that reflect real-world tool-calling scenarios: document processing, vehicle management, stock trading, travel planning, and file system management. These domains were selected to represent varying levels of complexity, different types of argument structures, and distinct sources of ambiguity that agents encounter in practice. Table 2 gives a statistical summary of the benchmark. Each sample in ClarifyBench is represented as a tuple: (*user query*, *user intent*, *follow-up queries*, *ground truth tool call*, *domain*).

The benchmark includes three distinct query types that systematically evaluate different aspects of clarification: **1. Explicit Queries:** Well-specified requests that provide sufficient information for direct tool execution, serving as baseline performance indicators. **2. Ambiguous Queries:** Requests with missing or unclear parameters that require clarification to determine the appropriate tool calls and arguments. **3. Infeasible Queries:** Requests which if executed at face value would generate errors due to invalid parameters, conflicting constraints, or impossible conditions.

### 4.2 BENCHMARK CONSTRUCTION

**Data Sources.** ClarifyBench draws from two primary sources to ensure diversity and realism. First, we extract successfully executed tool calls from the DocPilot Mathur et al. (2024), which provides real user interactions in document processing scenarios. Second, we leverage the Berkeley Function Calling Leaderboard (BFCL-v3) Patil et al. (2025), which offers data across multiple domains: vehicle control, stock trading, travel planning, and file system management.

Metric	Doc	Vehicle	Stocks	Travel	Files	All
Total Samples	181	139	143	119	134	716
Number of Tools	18	22	19	15	18	92
Avg # of Tool Calls	3.9	4.5	3.9	3.7	3.1	3.8
Explicit Queries	49	50	49	50	43	241
Ambiguous Queries	49	39	46	40	39	213
Infeasible Queries	48	49	38	18	45	198
Avg # of Follow-up	2.9	2.1	2.7	2.3	1.8	2.4

Table 2: Statistical summary ClarifyBench.

five alternative user queries that omit the obfuscated information. We also generate user intent prompts using in-context learning examples to capture the original tool call semantics. For infeasible queries, we design handwritten rules based on common API errors to create tool calls that would generate failures, followed by a similar LLM-based query augmentation process. We process *BFCL-v3* using existing explicit and ambiguous parameter queries from the benchmark, ensuring sample independence by removing cases with secondary API dependencies. We apply rule-based validation and LLM judgment (via in-context learning) to identify and exclude such cases. For retained samples, we strip secondary API utterances and tool calls from ground truth annotations. User intent prompts are generated through LLM processing, and infeasible queries are constructed using domain-specific rules, mirroring the DocPilot data strategy.

**Human Validation.** To ensure quality and naturalness, a human annotator evaluates all LLM-generated queries using three criteria: (A) naturalness of language, (B) faithfulness to expected tool calls including all required details while excluding obfuscated parameters, and (C) for infeasible queries, presence of error-inducing requirements. The annotator selects one optimal query per sample from the five generated alternatives.

## 5 STRUCTURED ARGUMENT UNCERTAINTY GUIDED ELICITATION AGENT

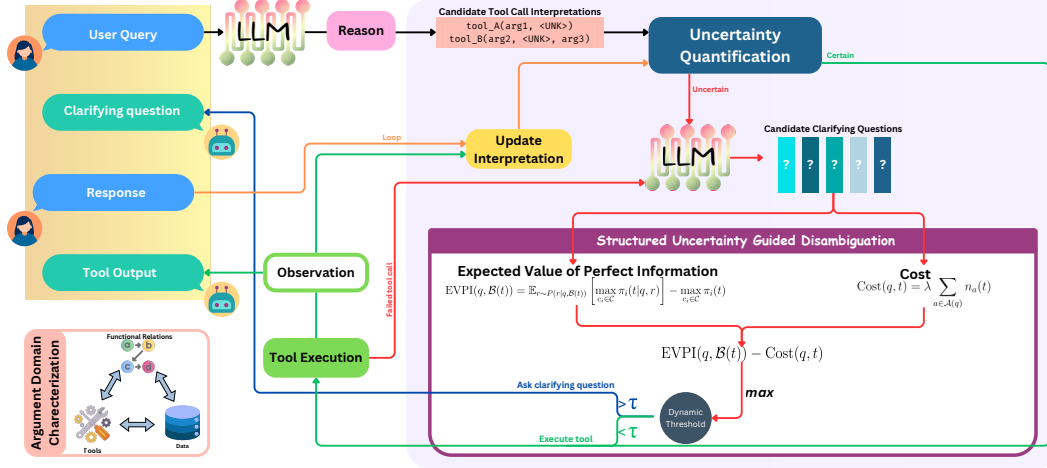


Figure 3: **SAGE-Agent**: ❶ Given a user query, an LLM reasons and generates potential tool calls with possibly uncertain parameters. These tool calls undergo ❷ structured uncertainty quantification to determine if clarification is needed. When uncertainty exists, the agent uses an LLM to produce ❸ candidate clarifying questions, and scores them using ❹ a cost-penalized Expected Value of Perfect Information (EVPI) metric. Tool-parameter domain interpretation is updated based on user-response to the clarifying question ❺, and given no further uncertainty, the best tool call is executed ❻.

SAGE (Structured Argument Uncertainty guided Elicitation) augments the standard Reason–Act–Observe loop by inserting structured, domain-aware clarification into the *Reason* stage (as seen in Fig. 3). Let the user input be  $u$ ; the toolkit  $\mathcal{T}$  and tool schemas follow Definition 1.

### 5.1 AGENT FLOW

At step  $t$ , the agent maintains belief  $\pi(t) = \{\pi_c(t)\}_{c \in \mathcal{C}}$  and observations  $\mathcal{O}_t$ . The full loop can be written as a combination of Reason ( $\mathcal{R}$ ) and Act ( $\mathcal{A}$ ):

$$(\mathcal{C}_t, Q_t) \xleftarrow{\mathcal{R}} (u, \mathcal{O}_t, \mathcal{T}) \xrightarrow{\mathcal{A}t} a_t = \begin{cases} \text{execute} : & c^*(t) = \arg \max_c \pi_c(t) \\ q^* : & \pi(t+1) = \mathcal{Ob}(\pi(t), o_{t+1}) \end{cases}$$

where  $\mathcal{R}$  produces candidate tool calls  $\mathcal{C}_t$  and aspect-targeted questions  $Q_t$ ,  $\mathcal{A}t$  selects either execution or query, and  $\mathcal{Ob}$  performs domain-constrained belief refinement (Fig. 3).

### 5.2 CANDIDATE GENERATION, QUESTIONING, AND BELIEF UPDATE

At step  $t$ , SAGE proceeds as follows:

**1. Candidate Generation.** The **Reason** stage prompts an LLM with  $(u, \mathcal{O}_t, \mathcal{T})$  to produce candidate tool calls  $\mathcal{C}_t = \{c_1, \dots, c_N\}$ , each assigning parameters  $\Theta_{i(c)}$  concrete values or <UNK>. Candidate certainty is defined as  $\tilde{\pi}_c(t) = \prod_{\theta_{i,j} \in \Theta_{i(c)}} p(\theta_{i,j} \mid T_{i(c)}, \text{obs}_t)$ . If  $\max_c \tilde{\pi}_c(t) \geq \tau_{\text{exec}}$ , execute  $c^*(t) = \arg \max_c \tilde{\pi}_c(t)$ ; otherwise continue.

**2. Question Generation.** An LLM is prompted with (i)  $q$ , (ii)  $\mathcal{C}$  and masks, (iii) tool schemas, and (iv) recent observations to output  $Q = \{(q_k, c_{i_k}, A_k)\}_{k=1}^L$ , where  $q_k$  is the question text,  $c_{i_k}$  the candidate being disambiguated, and  $A_k \subseteq \mathcal{A}$  the targeted aspects (parameters). Output is machine-parsable with <UNK> for ambiguous parameters.

**3. Scoring and Selection.** Let  $\mathcal{P}_q = \{C_1, \dots, C_M\}$  be the partition of  $\mathcal{C}_t$  induced by  $A$ . The EVPI is  $\text{EVPI}(q) = \sum_{m=1}^M \max_{c \in C_m} \tilde{\pi}_c(t) - \max_{c \in \mathcal{C}_t} \tilde{\pi}_c(t)$ . Score each question as  $\text{Score}(q, t) = \text{EVPI}(q) - \lambda \sum_{a \in A} n_a(t)$ , select  $q^*(t) = \arg \max_q \text{Score}(q, t)$ . If  $\max_q \text{Score}(q, t) < \alpha \max_c \tilde{\pi}_c(t)$  or budget  $n_s$  is exhausted, execute  $c^*(t)$ .

**4. Belief Update.** After observing answer  $r$ , update domains as  $\mathcal{D}_{i,j}(t+1) \leftarrow \mathcal{D}_{i,j}(t) \cap f_{\text{update}}(\theta_{i,j}, r)$  and recompute  $\tilde{\pi}_c(t+1)$ .

**5. Termination & Error Recovery.** Stop if (i)  $\max_c \tilde{\pi}_c(t) \geq \tau_{\text{exec}}$ , (ii)  $\max_q \text{Score}(q, t) < \alpha \max_c \tilde{\pi}_c(t)$ , or (iii)  $t \geq n_s$ . On execution failure, prompt for a fix or generate an error-specific  $q_{\text{error}}$  and re-enter step 3.

## 6 REWARD MODELING WITH STRUCTURED UNCERTAINTY

Our objective is to teach the agent not only *what* action to take but *when* to act with confidence versus request clarification. We fine-tune the policy using **Group Relative Policy Optimization (GRPO)** Shao et al. (2024), which samples multiple candidate actions per prompt, computes relative rewards, and updates the policy towards those exceeding the group mean—yielding a critic-free, memory-efficient variant of PPO that stabilizes optimization through implicit baselining and KL regularization. Our training data comes from the 9K examples in the When2Call Ross et al. (2025) dataset. For each user prompt and its tool set, the agent may take exactly one of four actions: AskQuestion, CallTool(parameters), Decline, or DirectAnswer. We prompt a base model to emit structured tags <reason>. . . </reason> <answer>. . . </answer>, and from that we compute scalar rewards.

### 6.1 BASELINE REWARD

The baseline reward is  $r_{\text{base}} = r_{\text{fmt}} + r_{\text{tool}} + r_{\text{cls}}$ , where  $r_{\text{fmt}} = 1.5$  (correct schema),  $r_{\text{tool}}$  equals 1.0 for correct tool+parameters, 0.75 if tool is correct but parameters are wrong, and 0.5 for correctly identifying a tool call or for non-tool actions, and  $r_{\text{cls}}$  equals up to 2.0 for correct action type. This encourages correctness and well-formedness but treats all instantiations equally regardless of model confidence or question informativeness.

### 6.2 CERTAINTY-WEIGHTED REWARD (OURS)

Let  $\pi_c(t)$  be the belief over candidate tool calls  $c \in \mathcal{C}_t$ . We define  $\text{Cert}(a_t) = \max_c \pi_c(t)$  if  $a_t$  is a tool call,  $1 - \max_c \pi_c(t)$  if  $a_t$  is a question, and 1 otherwise. The category reward becomes



$R_{\text{category}}(a_t) = \text{Cert}(a_t) \cdot r_{\text{base}}(a_t)$  which up-weights confident correct tool calls, penalizes low-certainty calls, and rewards clarification only when uncertainty is high—thus aligning reward with the agent’s own epistemic state.

🔑 **Key Insight:** Our reward is *self-calibrating*: it needs no critic to judge question quality, yet drives informative clarifications and confident tool calls. Unlike the baseline, which rewards all correct calls equally, our certainty-weighted reward **scales** with belief: confident calls get full payoff, low-confidence calls are penalized, and clarifications are rewarded only when uncertainty is high.

## 7 EXPERIMENTS

**ClarifyBench.** All baselines are implemented on a common ReAct agent scaffold for fair comparison. We evaluate **SAGE-Agent** against four baselines: (i) **ReAct + ask\_question()**, a standard ReAct agent with an ask\_question() tool serving as our control baseline; (ii) **ProCOT** Deng et al. (2023), which performs ProActive Chain-of-Thought reasoning to anticipate ambiguities before tool use; (iii) **Active Task Disambiguation** Kobalczyk et al. (2025), which generates candidate interpretations and clarification queries based on response entropy by parametrizing the solution space; and (iv) **Domain-aware ReAct**, which augments prompting and question generation with explicit schema information provided as context. All methods use GPT-4o and Qwen2.5-14B-Instruct with temperature 0.5. For SAGE-Agent, we pick  $\lambda = 0.5$ ,  $\alpha = 0.1$ ,  $\epsilon = 10^{-4}$ . We evaluate using four metrics: (1) **Coverage Rate**: proportion of tool calls with correct parameters matching the ground truth; (2) **Tool Match Rate (TMR)**: tool match rate against ground truth; (3) **Parameter Match Rate (PMR)**: parameter match rate against ground-truth; and (4) **Average Number of Questions (#Q)**: mean number of clarification questions asked per task (lower is better). We used 2xRTXA600 for inference.

**When2Call (Eval Set).** We trained GRPO with Qwen2.5-Instruct (3B and 7B) for one epoch using Unsloth. Three independent runs were performed, and results from the best-performing model are reported. Evaluation follows the original paper: log-probability comparison across options, option-prompted selection, and direct prompting without options. We trained using 4xL40S GPUs, and inferred using 1xL40S GPU. We train each setting for 3 runs, and report the setting with the best results.

## 8 RESULTS

Method	ClarifyBench - Ambiguous				ClarifyBench - Explicit				ClarifyBench - Infeasible			
	Coverage <sup>†</sup>	TMR <sup>†</sup>	PMR <sup>†</sup>	Avg #Q <sup>‡</sup>	Coverage <sup>†</sup>	TMR <sup>†</sup>	PMR <sup>†</sup>	Avg #Q <sup>‡</sup>	Coverage <sup>†</sup>	TMR <sup>†</sup>	PMR <sup>†</sup>	Avg #Q <sup>‡</sup>
<i>Base LLM: GPT-4o</i>												
ReAct + ask_question()	42.88 $\pm$ 25.1	70.41 $\pm$ 27.3	62.55 $\pm$ 23.9	2.68 $\pm$ 2.4	61.17 $\pm$ 22.7	87.95 $\pm$ 25.8	71.99 $\pm$ 28.4	2.15 $\pm$ 2.7	58.85 $\pm$ 24.3	85.05 $\pm$ 26.1	75.09 $\pm$ 21.8	2.21 $\pm$ 2.6
ProCOT	54.27 $\pm$ 27.4	75.62 $\pm$ 29.1	66.82 $\pm$ 24.6	2.07 $\pm$ 2.2	66.98 $\pm$ 22.8	89.57 $\pm$ 28.7	72.80 $\pm$ 25.4	2.14 $\pm$ 2.5	61.48 $\pm$ 24.2	89.32 $\pm$ 27.5	74.41 $\pm$ 23.5	2.43 $\pm$ 2.8
Active Task Disambiguation	45.60 $\pm$ 26.7	77.10 $\pm$ 28.2	60.78 $\pm$ 22.4	3.42 $\pm$ 2.6	66.97 $\pm$ 21.9	90.47 $\pm$ 29.3	72.45 $\pm$ 24.9	2.94 $\pm$ 2.5	65.27 $\pm$ 23.6	89.18 $\pm$ 28.8	75.09 $\pm$ 23.0	2.63 $\pm$ 2.3
Domain-aware ReAct	55.70 $\pm$ 24.5	79.83 $\pm$ 25.7	68.04 $\pm$ 23.3	2.56 $\pm$ 2.1	68.11 $\pm$ 22.5	91.17 $\pm$ 26.1	74.04 $\pm$ 25.2	2.10 $\pm$ 2.6	61.48 $\pm$ 24.0	90.32 $\pm$ 25.4	76.46 $\pm$ 26.7	2.03 $\pm$ 2.7
SAGE-Agent (Ours)	<b>59.73<math>\pm</math>22.1</b>	<b>86.02<math>\pm</math>27.5</b>	<b>71.79<math>\pm</math>25.3</b>	<b>1.39<math>\pm</math>2.0</b>	<b>71.67<math>\pm</math>21.8</b>	<b>93.65<math>\pm</math>29.7</b>	<b>75.94<math>\pm</math>26.1</b>	1.08 $\pm$ 2.2	<b>67.33<math>\pm</math>23.4</b>	<b>92.89<math>\pm</math>28.3</b>	<b>77.41<math>\pm</math>27.9</b>	<b>1.26<math>\pm</math>2.1</b>
<i>Base LLM: Qwen2.5-14B-Instruct</i>												
ReAct + ask_question()	40.34 $\pm$ 33.9	68.92 $\pm$ 32.0	63.35 $\pm$ 31.5	1.78 $\pm$ 1.94	51.85 $\pm$ 33.8	89.20 $\pm$ 22.8	73.63 $\pm$ 28.9	1.69 $\pm$ 1.67	42.39 $\pm$ 32.4	70.82 $\pm$ 31.1	63.31 $\pm$ 34.0	1.82 $\pm$ 1.43
ProCOT	52.45 $\pm$ 33.5	71.78 $\pm$ 33.7	70.08 $\pm$ 33.2	1.89 $\pm$ 2.03	61.76 $\pm$ 31.5	84.08 $\pm$ 23.8	74.60 $\pm$ 28.4	1.69 $\pm$ 1.68	52.08 $\pm$ 31.4	71.92 $\pm$ 29.3	68.72 $\pm$ 35.0	1.78 $\pm$ 1.51
Active Task Disambiguation	43.04 $\pm$ 29.2	69.06 $\pm$ 33.0	57.49 $\pm$ 34.1	2.45 $\pm$ 1.72	59.83 $\pm$ 33.1	81.01 $\pm$ 26.6	68.69 $\pm$ 31.5	2.31 $\pm$ 2.29	52.20 $\pm$ 30.6	76.59 $\pm$ 32.5	69.45 $\pm$ 35.0	2.22 $\pm$ 2.12
Domain-aware ReAct	51.10 $\pm$ 31.9	75.31 $\pm$ 30.7	67.50 $\pm$ 31.5	2.07 $\pm$ 1.35	60.91 $\pm$ 34.2	86.91 $\pm$ 24.8	71.70 $\pm$ 28.7	1.61 $\pm$ 1.56	55.76 $\pm$ 31.7	81.06 $\pm$ 27.2	72.23 $\pm$ 32.0	1.66 $\pm$ 1.30
SAGE-Agent (Ours)	<b>54.56<math>\pm</math>33.0</b>	<b>78.14<math>\pm</math>30.5</b>	<b>74.21<math>\pm</math>32.2</b>	<b>1.41<math>\pm</math>2.19</b>	<b>64.62<math>\pm</math>33.6</b>	<b>92.05<math>\pm</math>20.8</b>	<b>75.50<math>\pm</math>28.2</b>	<b>0.93<math>\pm</math>1.93</b>	<b>61.84<math>\pm</math>30.8</b>	<b>85.26<math>\pm</math>24.5</b>	<b>76.52<math>\pm</math>29.5</b>	<b>1.49<math>\pm</math>0.95</b>

Table 3: Performance comparison of agent strategies on ClarifyBench across two base LLMs (GPT-4o and Qwen2.5-14B-Instruct). **Best** and **second best** results within each LLM group are highlighted.

### 8.1 MAIN RESULTS: PERFORMANCE AND EFFICIENCY ON CLARIFYBENCH

Table 3 presents a comprehensive evaluation of SAGE-Agent against four baselines across three ClarifyBench categories (Ambiguous, Explicit, Infeasible) using two base LLMs (GPT-4o and



Qwen2.5-14B-Instruct). SAGE-Agent demonstrates consistent superiority across all evaluation dimensions, achieving state-of-the-art performance while simultaneously reducing user burden through fewer clarification questions.

**Performance Gains Across Task Categories.** On the Ambiguous split with GPT-4o, SAGE-Agent achieves 59.73% Coverage Rate, substantially outperforming the strongest baseline (Domain-aware ReAct at 55.70%). This 4.03 percentage point improvement extends to downstream metrics: Tool Match Rate (TMR) reaches 86.02% versus 79.83%, and Parameter Match Rate (PMR) attains 71.79% versus 68.04%. The pattern persists across Explicit scenarios, where SAGE-Agent achieves 71.67% Coverage (+3.56pp over Domain-aware ReAct), 93.65% TMR (+2.48pp), and 75.94% PMR (+1.90pp). Even on Infeasible tasks—where systems must recognize when queries cannot be satisfied—SAGE-Agent excels with 67.33% Coverage and 92.89% TMR, significantly outperforming all baselines. These gains demonstrate that structured schema-based reasoning enables more accurate task interpretation and parameter extraction than unstructured clarification approaches.

**Dramatic Reduction in User Burden.** A critical advantage of SAGE-Agent is its ability to achieve superior performance while asking dramatically fewer questions. On Ambiguous tasks with GPT-4o, SAGE-Agent averages just 1.39 questions per task—a 45.7% reduction compared to Domain-aware ReAct (2.56 questions) and an 48.1% reduction versus the basic ReAct baseline (2.68 questions). This reduction is even more pronounced compared to Active Task Disambiguation, which requires 3.42 questions per task on average. Critically, on Explicit scenarios where all necessary information is present in the initial query, SAGE-Agent asks only 1.08 questions—approaching the ideal of zero questions while still outperforming all baselines. The reduction in #Q directly translates to reduced user fatigue and improved user experience, as users provide fewer clarifications while receiving more accurate task execution.

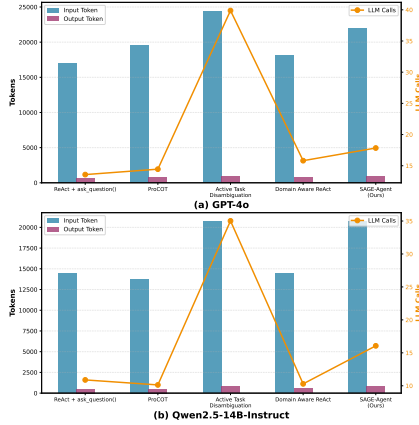


Figure 4: Resource consumption across methods for GPT-4o and Qwen2.5-14B.

**Computational Efficiency Despite Structured Reasoning.** Figure 4 reveals expected trade-offs: simpler baselines (ReAct, ProCOT, Domain-aware ReAct) exhibit lower footprints (14-18K input tokens, 14-16 calls with GPT-4o) but sacrifice performance (Table 3).

Among methods that explicitly model uncertainty, SAGE-Agent and Active Task Disambiguation represent distinct paradigms. Active Task Disambiguation quantifies uncertainty by generating candidate solutions and questions, then computing entropy over a  $|\text{questions}| \times |\text{solutions}|$  matrix—requiring 24K input tokens and 40 calls despite our token-efficient adaptation. SAGE-Agent instead parametrizes uncertainty directly over schema spaces using probabilistic reasoning, avoiding solution sampling entirely. This yields 22K input tokens with 54% fewer API calls, substantially reducing latency and cost while maintaining superior performance.

**Robustness Across Language Models.** SAGE-Agent’s advantages generalize across both proprietary and open-source LLMs. With Qwen2.5-14B-Instruct, SAGE-Agent achieves 54.56% Coverage on Ambiguous tasks, outperforming ProCOT (52.45%) and Domain-aware ReAct (51.10%), with TMR of 78.14% and PMR of 74.21%. The question reduction effect persists: 1.41 questions versus Domain-aware ReAct’s 2.07. On Explicit scenarios, SAGE-Agent reaches 64.62% Coverage with just 0.93 questions—a 42% reduction. On Infeasible tasks, SAGE-Agent achieves 61.84% Coverage and 85.26% TMR, substantially exceeding all baselines. While absolute metrics are lower with the smaller Qwen model compared to GPT-4o, the relative improvements over baselines remain consistent, demonstrating that structured schema-based clarification provides systematic advantages independent of base model choice.

**Impact of  $\lambda$ .** The redundancy penalty weight  $\lambda$  (Definition 5) controls the trade-off between information gathering and user burden by penalizing questions targeting previously queried aspects. Figure 5 shows the effect of  $\lambda \in \{0, 0.5, 1.0\}$  across 70 samples from each ClarifyBench split using GPT-4o, with independently scaled radar axes.

The results reveal a favorable operating point at  $\lambda = 0.5$ . Increasing  $\lambda$  from 0 to 0.5 yields substantial question reductions—18.1% on Ambiguous, 26.6% on Explicit, and 24.2% on Infeasible splits—while preserving task execution quality. Coverage Rate, TMR, and PMR remain stable with deviations under 3% across all settings, indicating that the penalized questions were indeed redundant rather than essential for task completion. The radar plots visualize this trade-off: the #Q dimension contracts inward while other metrics maintain consistent polygon shapes, demonstrating that question economy can be achieved without sacrificing accuracy.

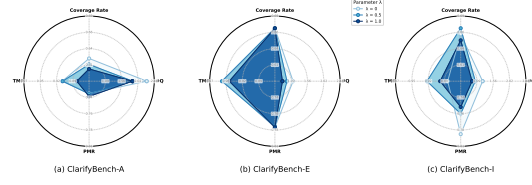


Figure 5: Effect of  $\lambda$  on performance metrics across ClarifyBench splits. Increasing  $\lambda$  from 0 to 0.5 reduces #Q by 18-27% while maintaining stable Coverage, TMR, and PMR ( $< 3\%$  deviation).

## 8.2 WHEN2CALL: LEARNING TO RECOGNIZE CLARIFICATION NEEDS

Figure 6 validates our hypothesis that uncertainty-aware training signals improve LLM clarification behavior. The When2Call benchmark tests models’ ability to recognize when clarification is needed versus when to proceed with available information.

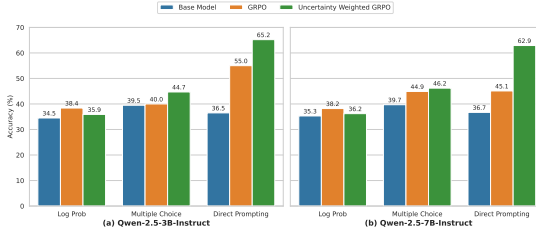


Figure 6: Performance of Qwen-2.5 models on When2Call across three evaluation methods: Log Probability, Multiple Choice, and Direct Prompting.

Comparing Qwen-2.5-3B and 7B models reveals that training signal quality matters more than model scale. The 3B model with uncertainty-weighted training (65.2% accuracy) substantially outperforms the 7B model with standard training (45.1% accuracy). This suggests that incorporating structured uncertainty into training objectives may be more valuable than simply scaling model parameters.

**Evaluation Mode Analysis.** The largest improvements occur in Direct Prompting mode, where models must make clarification decisions based solely on query analysis without multiple-choice scaffolding. This indicates that uncertainty-weighted training helps models develop robust internal representations of when clarification is needed, rather than merely improving selection among provided options.

## 9 CONCLUSION

Ambiguous user instructions fundamentally challenge tool-augmented LLM agents, leading to incorrect invocations and task failures. We presented **SAGE-Agent**, which models joint tool-argument clarification as a POMDP with Bayesian Value of Information objectives for optimal question selection. Extensive experiments validate our structured uncertainty approach: SAGE-Agent improves coverage on ambiguous tasks by 7–39% while reducing questions by 1.5–2.7 $\times$  on *ClarifyBench*, and uncertainty-weighted GRPO training boosts *When2Call* accuracy from 36.5% to 65.2% (3B) and 36.7% to 62.9% (7B). These results demonstrate that structured uncertainty provides a principled foundation for both inference and learning in tool-augmented scenarios. Our work establishes structured uncertainty quantification as essential for reliable, efficient LLM agents in real-world applications.

## 10 ETHICS STATEMENT

Our research does not use any personally identifiable information (PII) and all datasets employed in this work are used in accordance with their respective licenses (Apache 2.0). Our paper is designed primarily for deployment in collaborative AI assistance contexts where resolving ambiguity enhances productivity and user experience while minimizing unnecessary interaction. The system’s core approach of reducing clarification questions through principled uncertainty estimation promotes more equitable access to AI assistance by respecting users’ time and cognitive resources. While SAGE-Agent significantly reduces interaction burden, we recommend appropriate transparency about system limitations and human oversight when deploying in sensitive contexts. Furthermore, we encourage ongoing evaluation to ensure that question selection patterns do not reflect or amplify biases present in underlying models or training data. We acknowledge the ICLR code of ethics.

## REFERENCES

- Jayachandu Bandlamudi, Ritwik Chaudhuri, Neelamadhav Gantayat, Kushal Mukherjee, Prerna Agarwal, Renuka Sindhgatta, and Sameep Mehta. A framework for testing and adapting rest apis as llm tools. 2025. URL <https://arxiv.org/abs/2504.15546>.
- Michael Han Daniel Han and Unsloth team. Unsloth, 2023. URL <http://github.com/unslothai/unsloth>.
- Yang Deng, Wenqiang Lei, Wenxuan Zhang, Wai Lam, and Tat-Seng Chua. Pacific: towards proactive conversational question answering over tabular and textual data in finance. *arXiv preprint arXiv:2210.08817*, 2022.
- Yang Deng, Lizi Liao, Liang Chen, Hongru Wang, Wenqiang Lei, and Tat-Seng Chua. Prompting and evaluating large language models for proactive dialogues: Clarification, target-guided, and non-collaboration, 2023. URL <https://arxiv.org/abs/2305.13626>.
- Kung-Hsiang Huang, Akshara Prabhakar, Sidharth Dhawan, Yixin Mao, Huan Wang, Silvio Savarese, Caiming Xiong, Philippe Laban, and Chien-Sheng Wu. Crmarena: Understanding the capacity of llm agents to perform professional crm tasks in realistic environments. *arXiv preprint arXiv:2411.02305*, 2024a.
- Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024b.
- Katarzyna Kobalczyk, Nicolas Astorga, Tennison Liu, and Mihaela van der Schaar. Active task disambiguation with llms. *arXiv preprint arXiv:2502.04485*, 2025.
- Sizhe Liu, Yizhou Lu, Siyu Chen, Xiyang Hu, Jieyu Zhao, Yingzhou Lu, and Yue Zhao. Drugagent: Automating ai-aided drug discovery programming through llm multi-agent collaboration. *arXiv preprint arXiv:2411.15692*, 2024.
- Jiarui Lu, Thomas Holleis, Yizhe Zhang, Bernhard Aumayer, Feng Nan, Felix Bai, Shuang Ma, Shen Ma, Mengyu Li, Guoli Yin, et al. Toolsandbox: A stateful, conversational, interactive evaluation benchmark for llm tool use capabilities. *arXiv preprint arXiv:2408.04682*, 2024.
- Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn llm agents. *arXiv preprint arXiv:2401.13178*, 2024.
- Puneet Mathur, Alexa Siu, Varun Manjunatha, and Tong Sun. Docpilot: Copilot for automating pdf edit workflows in documents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pp. 232–246, 2024.
- Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.

- Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Zhong Zhang, Jie Zhou, Yankai Lin, Zhiyuan Liu, et al. Tell me more! towards implicit user intention understanding of language model driven agents. *arXiv preprint arXiv:2402.09205*, 2024.
- Tom Rainforth, Adam Foster, Desi R Ivanova, and Freddie Bickford Smith. Modern bayesian experimental design. *Statistical Science*, 39(1):100–114, 2024.
- Sudha Rao and Hal Daumé III. Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information. *arXiv preprint arXiv:1805.04655*, 2018.
- Hayley Ross, Ameya Sunil Mahabaleshwarkar, and Yoshi Suhara. When2Call: When (not) to call tools. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 3391–3409, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. URL <https://aclanthology.org/2025.naacl-long.174/>.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Sanidhya Vijayvargiya, Xuhui Zhou, Akhila Yerukola, Maarten Sap, and Graham Neubig. Interactive agents to overcome ambiguity in software engineering. 2025. URL <https://arxiv.org/abs/2502.13069>.
- Chenyu Wang, Weixin Luo, Sixun Dong, Xiaohua Xuan, Zhengxin Li, Lin Ma, and Shenghua Gao. 2025. URL <https://arxiv.org/abs/2401.10727>.
- Wenxuan Wang, Juluan Shi, Chaozheng Wang, Cheryl Lee, Youliang Yuan, Jen-tse Huang, and Michael R Lyu. Learning to ask: When llms meet unclear instruction. *arXiv preprint arXiv:2409.00557*, 2024.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. Travelplanner: A benchmark for real-world planning with language agents. *arXiv preprint arXiv:2402.01622*, 2024.
- Jingjing Xu, Yuechen Wang, Duyu Tang, Nan Duan, Pengcheng Yang, Qi Zeng, Ming Zhou, and Xu Sun. Asking clarification questions in knowledge-based question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 1618–1629, 2019.
- Guoli Yin, Haoping Bai, Shuang Ma, Feng Nan, Yanchao Sun, Zhaoyang Xu, Shen Ma, Jiarui Lu, Xiang Kong, Aonan Zhang, et al. Mmau: A holistic benchmark of agent capabilities across diverse domains. *arXiv preprint arXiv:2407.18961*, 2024.
- Yangyang Yu, Haohang Li, Zhi Chen, Yuechen Jiang, Yang Li, Denghui Zhang, Rong Liu, Jordan W Suchow, and Khaldoun Khashanah. Finmem: A performance-enhanced llm trading agent with layered memory and character design. In *Proceedings of the AAAI Symposium Series*, volume 3, pp. 595–597, 2024.
- Michael JQ Zhang and Eunsol Choi. Clarify when necessary: Resolving ambiguity through interaction with lms. *arXiv preprint arXiv:2311.09469*, 2023.
- Xuan Zhang, Yang Deng, Zifeng Ren, See Kiong Ng, and Tat-Seng Chua. Ask-before-plan: Proactive language agents for real-world planning. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 10836–10863, 2024.

## A METHOD DETAILS

### A.1 THEORETICAL PROOFS

**Proposition 1 (Viability Score Properties).** *The viability scoring function satisfies: (1) Monotonicity:  $\pi_i(t+1) \geq \pi_i(t)$  when information is gained, (2) Boundedness:  $0 \leq \pi_i(t) \leq 1$ , (3) Completeness:  $\pi_i(t) = 1$  iff all parameters are fully specified.*

*Proof.* (1) **Monotonicity:** Information gain can only constrain parameter domains:  $\mathcal{D}_{i,j}(t+1) \subseteq \mathcal{D}_{i,j}(t)$ . Therefore  $|\mathcal{D}_{i,j}(t+1)| \leq |\mathcal{D}_{i,j}(t)|$ , which implies  $|\mathcal{D}_{i,j}(t+1)|^{-1} \geq |\mathcal{D}_{i,j}(t)|^{-1}$ . Since  $\pi_i(t) = \prod_j p(\theta_{i,j})$  and each factor is non-decreasing,  $\pi_i(t+1) \geq \pi_i(t)$ .

(2) **Boundedness:** Each parameter certainty  $p(\theta_{i,j}) \leq 1$  by definition. Since  $\pi_i(t) = \prod_j p(\theta_{i,j})$ , we have  $0 \leq \pi_i(t) \leq 1$ .

(3) **Completeness:**  $\pi_i(t) = 1 \Leftrightarrow \prod_j p(\theta_{i,j}) = 1 \Leftrightarrow \forall j : p(\theta_{i,j}) = 1 \Leftrightarrow$  all parameters specified.  $\square$

**Proposition 2 (EVPI Properties).** *The EVPI function satisfies: (1) Non-negativity:  $\text{EVPI}(q, \mathcal{B}(t)) \geq 0$ , (2) Submodularity: diminishing returns for question sequences, (3) Convergence: EVPI approaches zero as uncertainty resolves.*

*Proof.* (1) **Non-negativity:** By Jensen’s inequality applied to the concave maximum function:

$$\mathbb{E}_r \left[ \max_{c_i} \pi_i(t|q, r) \right] \geq \max_{c_i} \mathbb{E}_r[\pi_i(t|q, r)] = \max_{c_i} \pi_i(t)$$

Therefore  $\text{EVPI}(q, \mathcal{B}(t)) \geq 0$ .

(2) **Submodularity:** For question sets  $S \subseteq S'$ , the marginal information gain satisfies:

$$\text{EVPI}(q|S) - \text{EVPI}(q|S') = H[\mathcal{B}|S] - H[\mathcal{B}|S \cup \{q\}] - (H[\mathcal{B}|S'] - H[\mathcal{B}|S' \cup \{q\}]) \geq 0$$

This follows from submodularity of entropy:  $H[X|Y] - H[X|Y, Z] \geq H[X|Y, W] - H[X|Y, W, Z]$  when  $W \supseteq \emptyset$ .

(3) **Convergence:** As uncertainty resolves,  $\max_i \pi_i(t) \rightarrow 1$  and candidate distributions become concentrated. For any question  $q$ ,  $\mathbb{E}_r[\max_i \pi_i(t|q, r)] \rightarrow \max_i \pi_i(t)$ , so  $\text{EVPI}(q) \rightarrow 0$ .  $\square$

**Theorem 1 (Finite Termination).** *Under regularity conditions on the response model, the algorithm terminates in finite expected time with probability 1.*

*Proof.* The termination condition is  $\max_q [\text{EVPI}(q) - \text{Cost}(q)] < \alpha \cdot \max_i \pi_i(t)$ .

**Case 1:** If  $\max_i \pi_i(t)$  increases over time (candidates improve), the right-hand side grows while EVPI values are bounded above. Eventually the inequality is satisfied.

**Case 2:** If  $\max_i \pi_i(t)$  remains bounded, then either: - EVPI values decrease due to information gain (Proposition 2.3) while costs increase linearly - Or no informative questions remain, making  $\text{EVPI} \approx 0$

In both cases, the net value becomes negative in finite time.

**Formal bound:** Let  $\rho = \mathbb{E}[\text{improvement in } \max_i \pi_i \text{ per question}]$  and  $\gamma = \mathbb{E}[\text{EVPI decline per question}]$ . - If  $\rho > 0$ : termination when  $\alpha \rho T \geq \text{EVPI}_{\text{initial}} - \gamma T$ , giving  $T \leq \frac{\text{EVPI}_{\text{initial}}}{\alpha \rho + \gamma}$  - If  $\rho \leq 0$ : termination when costs exceed EVPI, giving  $T \leq \frac{\max \text{EVPI}}{\lambda \cdot \min |\mathcal{A}(q)|}$

Therefore  $\mathbb{E}[T] < \infty$ .  $\square$

### A.2 SAGE-AGENT: FURTHER DETAILS

The following algorithm presents the complete practical implementation of our Structured Agent Guided Elicitation (SAGE) framework, which operationalizes the theoretical principles developed in Section 2:

**Algorithm 1** SAGE-Agent**Require:** user input  $u$ , toolkit  $\mathcal{T}$ , max questions  $n_s$ , redundancy weight  $\lambda$ , thresholds  $\tau_{\text{exec}}, \alpha$ 

```

1:  $(\mathcal{C}_0, \mathcal{O}_0) \leftarrow \mathcal{R}(u, \emptyset, \mathcal{T})$ 
2: compute  $w_c(0)$  for  $c \in \mathcal{C}_0$ ; normalize  $\pi_c(0) = w_c(0) / \sum_{c'} w_{c'}(0)$ 
3:  $t \leftarrow 0$ 
4: while true do
5:   if  $\max_c \pi_c(t) \geq \tau_{\text{exec}}$  then
6:     execute  $c^* = \arg \max_c \pi_c(t)$ ; break
7:   end if
8:    $Q_t \leftarrow \text{LLM\_GenerateQuestions}(u, \mathcal{C}_t, \mathcal{O}_t)$ 
9:   for all  $q \in Q_t$  do
10:    compute partition  $\mathcal{P}_q$  over  $\mathcal{C}_t$  induced by  $\mathcal{A}(q)$ 
11:    compute  $\text{EVPI}(q)$ 
12:    compute  $\text{Cost}(q, t)$ 
13:     $\text{Score}(q, t) \leftarrow \text{EVPI}(q) - \text{Cost}(q, t)$ 
14:   end for
15:    $q^* \leftarrow \arg \max_{q \in Q_t} \text{Score}(q, t)$ 
16:   if  $\text{Score}(q^*, t) < \alpha \cdot \max_c \pi_c(t)$  or  $t \geq n_s$  then
17:     execute  $c^* = \arg \max_c \pi_c(t)$ ; break
18:   end if
19:   ask  $q^*$ ; observe response  $r$ ; update  $\mathcal{O}_{t+1} \leftarrow \mathcal{O}_t \cup \{r\}$ 
20:   increment  $n_a(t)$  for  $a \in \mathcal{A}(q^*)$ 
21:   refine affected domains:  $\mathcal{D}_{i,j}(t+1) \leftarrow \mathcal{D}_{i,j}(t) \cap f_{\text{update}}(\theta_{i,j}, r)$ 
22:   recompute  $w_c(t+1)$  and normalize to  $\pi_c(t+1)$ 
23:   if execution of current best later fails then
24:     generate  $q_{\text{error}} \leftarrow f_{\text{error}}(\cdot)$  and treat it like other  $q$ 
25:   end if
26:    $t \leftarrow t + 1$ 
27: end while

```

## A.2.1 CRITICAL IMPLEMENTATION DETAILS

**Domain Constraint Propagation.** The domain refinement step implements the constraint extraction function  $f_{\text{update}}(\theta_{i,j}, r)$  that maps natural language responses to parameter domain constraints. This function must handle:

- **Explicit constraints:** Direct specifications like "departure date is March 15th"
- **Schema Dependency:** When value of a specific parameter constrains available options for another parameter. Could be computed based on data.
- **Negative constraints:** Exclusions like "not business class"  $\rightarrow \text{class} \in \{\text{economy}, \text{premium}\}$

**Error Recovery Mechanism.** When the highest-confidence candidate fails at runtime, the system generates diagnostic questions  $q_{\text{error}}$  using function  $f_{\text{error}}(\cdot)$ . This adaptive questioning strategy enables recovery from:

- API failures or timeouts
- Invalid parameter combinations that pass initial validation

## A.2.2 HYPERPARAMETER SELECTION GUIDELINES

**Cost Parameter  $\lambda$ :** Controls trade-off between information gain and question burden - Small  $\lambda$  ( $< 0.1$ ): Aggressive questioning, may annoy users - Large  $\lambda$  ( $> 1.0$ ): Conservative questioning, may under-clarify - Recommended:  $\lambda \in [0.2, 0.5]$  based on empirical evaluation

**Termination Threshold  $\alpha$ :** Controls when to stop asking questions - Small  $\alpha$  ( $< 0.1$ ): Early termination, may execute with uncertainty - Large  $\alpha$  ( $> 0.5$ ): Late termination, more questions asked - Recommended:  $\alpha \in [0.1, 0.3]$  depending on task criticality

**Algorithm 2** ClarifyBench Agent Simulation Scaffold: Multi-Request Simulation Process

---

```

1: procedure EXECUTESIMULATION( $\mathcal{S}$ ) ▷  $\mathcal{S}$  represents the simulation scenario
2:   Initialize agent  $\mathcal{A}$ , environment  $\mathcal{E}$ , user model  $\mathcal{U}$ 
3:    $\mathcal{R} \leftarrow \{r_0, r_1, \dots, r_n\}$  ▷ Request sequence
4:    $\mathcal{C} \leftarrow \emptyset$  ▷ Conversation history
5:   for each request  $r_i \in \mathcal{R}$  do
6:      $\mathcal{T}_i \leftarrow \emptyset$  ▷ Turn sequence for request  $i$ 
7:      $q_{\text{current}} \leftarrow r_i$  ▷ Current query state
8:      $\text{clarification\_count} \leftarrow 0$ 
9:     while  $\text{clarification\_count} < \tau_{\text{max}}$  and not terminated do
10:       $\text{response} \leftarrow \mathcal{A}(q_{\text{current}}, \mathcal{C})$ 
11:      if  $\text{response} \in \Phi_{\text{success}}$  then ▷ Successful completion
12:        Record completion in  $\mathcal{T}_i$ 
13:        break
14:      else if  $\text{response} \in \Phi_{\text{clarification}}$  then ▷ Needs clarification
15:         $\text{clarification} \leftarrow \mathcal{U}(\text{response.question}, \mathcal{S})$ 
16:        if  $\text{clarification} = \perp$  then ▷ User cannot provide clarification
17:          Record incomplete in  $\mathcal{T}_i$ 
18:          break
19:        end if
20:         $q_{\text{current}} \leftarrow \text{Enrich}(r_i, \text{clarification})$ 
21:         $\text{clarification\_count} \leftarrow \text{clarification\_count} + 1$ 
22:      else
23:        Record failure in  $\mathcal{T}_i$ 
24:        break
25:      end if
26:    end while
27:     $\mathcal{C} \leftarrow \mathcal{C} \cup \mathcal{T}_i$ 
28:  end for
29:   $\text{evaluation} \leftarrow \text{Evaluate}(\mathcal{C}, \mathcal{S}.\text{ground\_truth})$ 
30:  return  $\{\mathcal{C}, \text{evaluation}\}$ 
31: end procedure

```

---



### A.3 GROUP RELATIVE POLICY OPTIMIZATION (GRPO) EXPERIMENT

#### A.3.1 DATASET CREATION AND PROCESSING

**Source Dataset:** Our enhanced dataset was constructed from the nvidia/When2Call dataset, from the "train\_pref" data. This dataset contains preference-ranked examples for tool-calling tasks with human-annotated preferred responses for training reinforcement learning models.

**Original Data Structure:** Each example in the source dataset contained:

- **Messages:** Conversation history with user and assistant exchanges in chat format
- **Tools:** Available tool definitions with JSON schema parameters and descriptions
- **Chosen responses:** Human-preferred responses for the given context
- **Preference annotations:** Quality ratings for different response options

**Response Classification:** Each example was processed to classify responses into four categories: <TOOLCALL>, <ASK>, <REFUSE>, and <DIRECTLY>. Classification used keyword-based heuristics:

- <TOOLCALL>: Presence of "<TOOLCALL>" tags or "toolcall" keywords
- <ASK>: Presence of question marks ("?",) in content
- <REFUSE>: Presence of refusal keywords ("sorry", "unable", "impossible", etc.)
- <DIRECTLY>: Default classification for other responses. (None existed in the preferred set)

**Data Transformations:** Several preprocessing steps were applied to optimize the dataset for uncertainty-aware training:

1. **Domain Schema Injection:** Each example was augmented with parsed domain information for all available tools, stored as JSON strings in a tool\_domain\_schemas field for HuggingFace compatibility
2. **Message Format Preservation:** The chat format was maintained with modified system messages while preserving user/assistant alternation

#### A.3.2 TOOL DOMAIN ANALYSIS

To enable uncertainty quantification, we performed comprehensive domain analysis of all available tools using Qwen-2.5-7B-Instruct as the primary analysis model. Each tool's arguments were analyzed to determine:

- **Domain type:** finite, estimated\_finite, numeric\_range, string, boolean, list, or custom
- **Domain size:** exact count for finite domains, estimates for larger domains, or infinite for unbounded domains
- **Domain values:** complete enumeration for small domains, representative examples for larger domains, or range bounds for numeric domains
- **Data dependency:** whether argument values depend on external data sources or user context

The analysis prompt instructed the model to classify arguments according to strict validation rules:

- Finite domains ( $\leq 20$  values): complete value enumeration with domain\_size = len(domain\_values)
- Estimated finite domains: 5-10 representative examples with domain\_size  $\gg$  len(examples)
- Numeric ranges: [min, max] bounds with appropriate size calculation
- Boolean domains: domain\_size = 2 with null values
- String/custom domains: infinite size with null values

### A.3.3 UNCERTAINTY-AWARE SYSTEM PROMPTS

Each training example was enhanced with a comprehensive system prompt that provided explicit instructions for uncertainty handling. The complete system prompt template was:

```
\texttt{You are a helpful agent. You will have access to tools to answer
the query.\\
\\
UNCERTAINTY GUIDELINES:\\
- Use <UNK> for arguments you cannot determine from context, or cannot
reasonably estimate. Don't overuse, you can assume defaults where
needed.\\
- When asking questions, use the structured format with candidate tool
calls\\
\\
You can perform following action types:\\
a) <TOOLCALL> Invoke a tool call as follows:\\
<TOOLCALL>\\
[\\{"name": "tool\\_name", "arguments": \\{"argument\\_name": "value", "
uncertain\\_argument": "<UNK>", ...\\}\\}\\]\\
</TOOLCALL>\\
\\
b) <ASK> Ask a question from the user if you need more information to
execute a tool call </ASK>\\
\\
STRUCTURED QUESTION FORMAT (when asking for clarification):\\
<ASK>\\
<TOOLCALL>\\
// Think about what tool you would call given the request, and the
current information. Because some information is missing, you want to
ask a question.\\
[
\\{\\{ "name": "tool\\_name", "arguments": \\{"known\\_arg": "value", "
uncertain\\_arg": "<UNK>"\\}\\}\\]\\
</TOOLCALL>\\
<question>\\
What is the specific value for uncertain\\_arg?\\
</question>\\
</ASK>\\
\\
c) <REFUSE> Refuse, if your knowledge or available tools can't be used
here </REFUSE>\\
d) <DIRECTLY> directly answer </DIRECTLY>\\
\\
Your response should be formatted like:\\
<reasoning>\\
Step-by-step thinking about certainty/uncertainty of each argument\\
</reasoning>\\
<answer>\\
<ACTION\\_TYPE>\\
..content.. (Question/ToolCall/Refuse/DirectAnswer)\\
</ACTION\\_TYPE>\\
</answer>}
```

### A.3.4 TRAINING CONFIGURATION

Training began from unsloth/Qwen2.5-3B-Instruct and unsloth/Qwen2.5-7B-Instruct checkpoints. LoRA (Low-Rank Adaptation) fine-tuning was applied with rank 64 adaptations targeting attention and MLP projection layers.

Model training was performed using Group Relative Policy Optimization, using Unsloth Daniel Han & team (2023) with parameter details in Table 4.

Hyperparameter	Value
Learning Rate	5e-6
Per Device Batch Size	1 (3B), 8 (logs)
Gradient Accumulation Steps	1
Max Sequence Length	1024
Training Epochs	1
Warmup Ratio	0.1
Weight Decay	0.1
Optimizer	AdamW 8-bit
Adam Beta1	0.9
Adam Beta2	0.99
LoRA Rank	64
LoRA Alpha	64

Table 4: Training hyperparameters for uncertainty-aware tool calling model.

### A.3.5 REWARD MODELING

Our baseline GRPO reward function consists of multiple components that guide the model toward generating well-formed, accurate responses. The total reward for a generated completion is computed as the sum of three independent reward components:

$$r_{\text{total}} = r_{\text{fmt}} + r_{\text{tool}} + r_{\text{cls}} \quad (7)$$

where  $r_{\text{fmt}}$  represents format compliance rewards,  $r_{\text{tool}}$  represents tool call accuracy, and  $r_{\text{cls}}$  represents action classification rewards.

**Format Compliance Rewards ( $r_{\text{fmt}}$ ).** These components encourage proper XML formatting and total up to 1.5 points:

- **XML Count Reward:** Awards up to 0.5 points for proper newline structure, penalizing excessive trailing content.
- **Soft Format Reward:** Awards 0.5 points if the response contains `<reasoning>` and `<answer>` tags in the correct order (with flexible whitespace).
- **Strict Format Reward:** Awards 0.5 points only if the response exactly matches the format `<reasoning>\n...\n</reasoning>\n<answer>\n...\n</answer>\n`.

**Tool Call Accuracy Reward ( $r_{\text{tool}}$ ).** Compares the predicted tool call against a ground truth reference:

$$r_{\text{tool}} = \begin{cases} 1.0 & \text{if tool name and arguments match exactly} \\ 0.75 & \text{if tool name matches but arguments differ} \\ 0.5 & \text{if both have no tool call OR wrong tool name} \\ 0.0 & \text{if one has a tool call and the other does not} \end{cases} \quad (8)$$

**Action Classification Reward ( $r_{\text{cls}}$ ).** This reward is the primary component that differentiates between GRPO and Certainty weighted GRPO. This reward is computed based on the agent’s chosen action  $a_t$  at timestep  $t$ , which can be: TOOLCALL (execute a tool), ASK (request clarification), REFUSE (decline the request), or DIRECTLY (answer without tools).

The base classification reward is computed as:

$$r_{\text{cls}}(a_t) = \begin{cases} 2.0 & \text{if response starts with correct tag and contains } \geq 30 \text{ chars} \\ 1.5 & \text{if response starts with correct tag but insufficient content} \\ 0.0 & \text{otherwise} \end{cases} \quad (9)$$

**Certainty Weighting** For the baseline **GRPO**, the final classification reward is simply:

$$r_{\text{cls}}^{\text{GRPO}}(a_t) = r_{\text{cls}}(a_t) \quad (10)$$

For **Certainty weighted GRPO**, we introduce epistemic-state-aware weighting. Let  $\pi_c(t)$  be the model’s belief over candidate tool calls  $c \in \mathcal{C}_t$ . We define the certainty function:

$$\text{Cert}(a_t) = \begin{cases} \max_c \pi_c(t) & \text{if } a_t \text{ is a tool call} \\ 1 - \max_c \pi_c(t) & \text{if } a_t \text{ is a clarification question} \\ 1 & \text{otherwise} \end{cases} \quad (11)$$

The final classification reward is then:

$$r_{\text{cls}}^{\text{Certainty}}(a_t) = \text{Cert}(a_t) \cdot r_{\text{cls}}(a_t) \quad (12)$$

This formulation up-weights confident correct tool calls, penalizes low-certainty calls, and rewards clarification only when uncertainty is high—thus aligning the reward with the agent’s own epistemic state.

In our implementation, we approximate  $\pi_c(t)$  through explicit certainty computation over tool call arguments. For a tool call  $c$  with arguments, the certainty is:

$$\pi_c(t) = \prod_{\text{arg} \in c.\text{arguments}} \pi_{\text{arg}} \quad (13)$$

where for each argument:

$$\pi_{\text{arg}} = \begin{cases} 1.0 & \text{if arg has a specified value} \\ \frac{1}{|\mathcal{D}_{\text{arg}}|} & \text{if arg is empty and domain size is finite} \\ \epsilon \approx 0.0001 & \text{if arg is empty and domain size is infinite} \end{cases} \quad (14)$$

Here,  $\mathcal{D}_{\text{arg}}$  represents the domain size for that argument as specified in the tool schema. This approach ensures that tool calls with all arguments specified receive maximum certainty ( $\pi_c(t) = 1.0$ ), while tool calls with missing arguments receive certainty inversely proportional to the domain sizes of unspecified parameters. For ASK actions, we compute certainty over the candidate tool call mentioned in the question, and use  $1 - \pi_c(t)$  to reward asking when uncertainty is high.

#### A.4 PROMPT TEMPLATES

##### A.4.1 REACT AGENT PROMPTS

**Reasoning Prompt** This prompt is used in the main reasoning phase of the ReAct agent to decide which tool to use next based on the current state of the conversation.

```

You are an AI assistant helping with a user request.
SYSTEM CONTEXT:
You have access to the following tool domain:
{plugin_descriptions}
Request: {request}
Previous observations:
{obs_text}
Available tools:
{tool_registry.get_tool_descriptions()}
Think step by step about what tool to use next. Consider the plugin
context above to understand the capabilities available to you. If you
have enough information to provide a final answer, use the
final_answer tool.
Respond in JSON format:
{
  "reasoning": "Your step-by-step thinking",
  "tool_call": {
    "tool_name": "name_of_tool",
    "arguments": {
      "arg1": "value1",
      "arg2": "value2"
    }
  }
}

```

**Error Recovery Prompt** Used when a tool execution fails to determine if the error can be resolved automatically.

```
You are helping fix a failed tool call.
Original Request: {request}
Tool Information:
{tool_info or f"Tool: {tool_name}"}
Error Details:
{error_result.message}
Based on the error and tool information, can you suggest how to fix this?
Respond in JSON format:
{
  "can_fix": true/false,
  "reasoning": "explanation of what went wrong and how to fix it",
  "suggested_action": "retry_with_changes" or "different_tool" or "
    need_clarification",
  "observation": "observation to add to context for next reasoning step"
}
If you cannot determine a fix from the available information, set can_fix
to false.
```

#### A.4.2 QUESTION GENERATION PROMPT

Used to generate clarification questions when there is uncertainty about tool arguments.

```
You are an AI assistant that helps users by understanding their queries
and executing tool calls.
{conversation_history}Original user query:
"{user_query}"
Based on the query, I've determined that the following tool calls are
needed, but some arguments are uncertain:
Tool Calls:
{tool_calls}
Detailed Tool Documentation:
{tool_documentation}
Uncertain Arguments:
{uncertain_args}
Your task is to generate clarification questions that would help resolve
the uncertainty about specific arguments.
Instructions:

Generate questions that are clear, specific, and directly address the
uncertain arguments
Each question should target one or more specific arguments
Questions should be conversational and easy for a user to understand
For each question, specify which tool and argument(s) it aims to clarify.
Generate 5 diverse questions.
Keep in mind the the arguments you wish to clarify, their domains etc.

Return your response as a JSON object with the following structure:
{
  "questions": [
    {
      "question": "A clear question to ask the user",
      "target_args": [{"tool_name", "arg_name"}, {"tool_name", "other_arg_name"}]
    }
  ]
  // ... 5 total questions
}
Ensure that each question targets at least one uncertain argument.
```

### A.4.3 USER SIMULATOR

The simulator takes a language model provider, ground truth data, and user intent as inputs. It maintains the conversation state and ensures responses are consistent with the user's information. The core of the simulation lies in two prompt templates that instruct a language model to act as a user:

```
You are simulating a user who is interacting with an AI assistant.
Original query: "{self.original_query}"
User's intent for the CURRENT request: {self.user_intent}
Information needed for the CURRENT request (do not reveal future
intentions):
{current_turn_ground_truth}
Additional context:
{self.context}
The AI assistant has asked the following specific question:
"{question}"
Generate a realistic user response to this SPECIFIC question. The
response should:

Be natural and conversational
ONLY provide information that directly answers the specific question
asked
NOT mention any future requests or intentions the user might have
ONLY focus on the current task, not on future tasks
Be concise and to the point

IMPORTANT: Never reveal future intentions. Respond ONLY to the specific
question asked.
NEVER BREAK CHARACTER. DO NOT THINK OUT LOUD. Respond directly as the
user would:
```

This template ensures the simulator provides natural, conversational responses that only address the specific question without revealing future intentions. For generating follow-up requests, the simulator uses this template:

```
You are simulating a user who is interacting with an AI assistant.
Original query: "{self.original_query}"
User's intent: {self.user_intent}
Previous conversation:
{formatted_history}
Based on the conversation so far and the user's intent, decide if the
user would have a follow-up request.
Consider:

Has everything the user wanted been accomplished?
Is there a logical next step the user might want to take?
Has the agent clearly indicated that they've completed all necessary
tasks?

If you believe the user would have a follow-up request, provide it in a
natural, conversational way.
If you believe the conversation is complete, respond with "
CONVERSATION_COMPLETE".
NEVER BREAK CHARACTER, DO NOT THINK!
Decision:
```

This template helps the simulator determine whether to generate a follow-up request based on the conversation context and predefined potential follow-ups. The User Simulator isolates ground truth information for each conversation turn, ensuring only relevant information is revealed at appropriate times. It tracks the original query, user intent, ground truth for tool calls, completed tool calls, potential follow-up queries, and the current conversation turn. By providing consistent, realistic user responses, the simulator allows for reproducible evaluation of clarification strategies across multiple scenarios.

## B BENCHMARK DETAILS

### B.1 BENCHMARK DOMAINS

This appendix describes the key characteristics of each API domain used in our experiments, detailing their initialization parameters, state management, and tool specifications.

**Gorilla File System Plugin (GFS).** The Gorilla File System API simulates a UNIX-like file system with a hierarchical directory structure. It maintains state through:

- Directory structure with nested files and subdirectories
- Current working directory pointer
- Each file contains content as strings

The plugin provides 18 tools implementing common file system operations such as navigation, file creation, modification, and content manipulation. Each tool supports parameters relevant to file system operations, such as file names, directory paths, and content strings. Table 9 provides detailed information about these tools and their parameter domains.

The GFS plugin’s domains depend heavily on the current state of the file system. Domain updates revolve primarily around available files and directories in the current working directory, as outlined in Table 10.

**Document Processing.** The Document API simulates operations for PDF document manipulation. Its state consists of:

- Number of pages in the current document
- PDF filename metadata
- Operation-specific context for page-based operations

The plugin provides 18 document manipulation tools including conversion, annotation, redaction, and page manipulation functions. Parameters include page numbers, text content, formatting options, and file paths. Table 6 details the tools and their parameter domains.

Domain updates in the Document Plugin focus on page numbers and ranges, adapting dynamically to changes in document length when pages are added or deleted, as shown in Table 10.

**Vehicle Control.** The Vehicle Control API simulates an automotive control system with:

- Engine state (running or stopped)
- Door lock status for each door
- Fuel level (ranging from 0 to 50 gallons)
- Battery voltage
- Climate control settings
- Brake systems (pedal position and parking brake)
- Lighting systems
- Navigation state

This plugin implements 24 vehicle control tools that manipulate different aspects of the vehicle, including engine operations, door management, climate control, lighting, braking systems, and navigation. Table 8 details the specific tools and their parameter domains.

Vehicle Control domain updates primarily concern contextual constraints such as brake pedal position for engine start, door states, and fuel level requirements, as referenced in Table 10.



**Travel.** The Travel API simulates a travel booking and management system with:

- Credit card registry and balances
- Flight booking records
- User information (first name, last name)
- Budget limits
- Available routes with pricing data

The plugin provides 15 tools for travel-related operations, including flight bookings, credit card management, budget settings, and travel information queries. Table ?? details these tools and their parameter domains.

Domain updates in the Travel Plugin focus on available credit cards, booking IDs, and airport codes for valid routes, as detailed in Table 10.

**Trading Bot.** The Trading Bot simulates a stock trading platform with:

- Account information and balance
- Order records (pending, completed, cancelled)
- Stock data with prices and metrics
- Watchlist of stocks
- Transaction history
- Market status (open/closed)

This plugin provides 19 trading tools for account management, order placement, stock information retrieval, and market analysis. Table 7 lists the specific tools and their parameter domains.

Trading Plugin domain updates primarily involve available stocks, watchlist items, and order IDs, adapting to user actions like placing orders or modifying watchlists, as referenced in Table 10.

All plugins follow a consistent pattern for state initialization through configuration objects, domain updates based on state changes, and parameter validation. The dynamic nature of these domains presents particular challenges for language model interactions, as valid parameter values continuously evolve during conversations based on system state changes.

## B.2 HUMAN ANNOTATION

We employed two graduate student annotators, aged 22-25. The annotators were proficient in English, and have proficiency in Python (relevant to test tool calls). The annotators were fairly compensated at the standard Graduate Assistant hourly rate, following their respective graduate school policies. Fig 7 shows a summary of the annotator guidelines.

## B.3 TOOL CALL CORRUPTION HEURISTICS

We handcrafted rules to corrupt validated tool calls in the ground truth data, to construct ClarifyBench-Infeasible.

**GorillaFileSystem** For the file system API, we implemented four primary corruption strategies:

- *Invalid File Name Corruption* targeting functions like `mkdir`, `touch`, and `cat` by inserting forbidden characters (e.g., `|`, `/`, `\`, `?`);
- *Path Traversal Corruption* for `cd`, `mv`, `cp`, and `find` operations by inserting relative paths (`./`, `../`) or absolute paths (`/root/`);
- *Non-existent Files Corruption* for file operation functions by generating random names or modifying existing names;
- *Duplicate Creation Corruption* for `mkdir` and `touch` operations by using existing file/directory names.

Tool Name	Argument	Description	Domain Type	Domain Values	Data Dep.	Required
get_budget_fiscal_year	lastModifiedAfter	Date filter for fiscal years	string	Any date string	N	N
	includeRemoved	Include removed fiscal years	string	Any string	N	N
register_credit_card	card_number	Credit card number	string	Any card number	N	Y
	expiration_date	Card expiration (MM/YYYY)	string	MM/YYYY format	N	Y
	cardholder_name	Name on card	string	Any name string	N	Y
	card_verification_number	CVV code	numeric_range	[100, 999]	N	Y
get_flight_cost	travel_from	Departure airport code	string*	3-letter codes	Y	Y
	travel_to	Arrival airport code	string*	3-letter codes	Y	Y
	travel_date	Travel date	string	YYYY-MM-DD	N	Y
	travel_class	Seat class	finite	[economy, business, first]	N	Y
get_credit_card_balance	card_id	Credit card identifier	string*	Card ID list	Y	Y
book_flight	card_id	Payment card ID	string*	Card ID list	Y	Y
	travel_date	Travel date	string	YYYY-MM-DD	N	Y
	travel_from	Departure airport	string*	Airport codes	Y	Y
	travel_to	Arrival airport	string*	Airport codes	Y	Y
	travel_class	Seat class	finite	[economy, business, first]	N	Y
	travel_cost	Flight cost	numeric_range	[0, 10000]	N	Y
retrieve_invoice	booking_id	Booking identifier	string*	Booking ID list	Y	N
	insurance_id	Insurance identifier	string*	Insurance ID list	Y	N
list_all_airports	No arguments					
cancel_booking	booking_id	Booking to cancel	string*	Booking ID list	Y	Y
compute_exchange_rate	base_currency	Source currency	finite	[USD, RMB, EUR, JPY, GBP, CAD, AUD, INR, RUB, BRL, MXN]	N	Y
	target_currency	Target currency	finite	[USD, RMB, EUR, JPY, GBP, CAD, AUD, INR, RUB, BRL, MXN]	N	Y
	value	Amount to convert	numeric_range	[0, 1000000]	N	Y
verify_traveler_information	first_name	Traveler's first name	string	Any name	N	Y
	last_name	Traveler's last name	string	Any name	N	Y
	date_of_birth	Birth date	string	YYYY-MM-DD	N	Y
	passport_number	Passport number	string	Any passport ID	N	Y
set_budget_limit	budget_limit	Budget limit in USD	numeric_range	[0, 10000]	N	Y
get_nearest_airport_by_city	location	City name	finite	[Rivermist, Stonebrook, ...]	N	Y
purchase_insurance	insurance_type	Type of insurance	finite	[basic, premium, deluxe]	N	Y
	booking_id	Booking identifier	string*	Booking ID list	Y	Y
	insurance_cost	Insurance cost	numeric_range	[0, 1000]	N	Y
	card_id	Payment card ID	string*	Card ID list	Y	Y
contact_customer_support	booking_id	Booking reference	string*	Booking ID list	Y	Y
	message	Support message	string	Any message text	N	Y
get_all_credit_cards	No arguments					

Table 5: Travel Plugin API: Complete Tool and Argument Specification with Domain Dependencies (without Importance column)

**DocumentPlugin** For the document manipulation API, we implemented three corruption strategies:

- *Invalid Page Range Corruption* for functions like `add_comment` and `delete_page` by setting zero/negative values or exceeding total pages;
- *Invalid Formats Corruption* for convert operations by using unsupported formats or partial strings;
- *Out of Range Values Corruption* for parameters like `font_size` and `transparency` by exceeding min/max bounds or using negative values.

**VehicleControlAPI** For the vehicle control API, we focused on two corruption categories:

- *Invalid Ranges Corruption* for functions like `fillFuelTank` and `adjustClimateControl` by exceeding capacity or using negative values;
- *Invalid Enums Corruption* for operations like `startEngine` and `setHeadlights` by supplying wrong enum values or case mismatches.

**TravelAPI** For the travel booking API, we implemented three corruption strategies:

- *Financial Constraints Corruption* for functions like `book_flight` by exceeding available balance or using negative values;
- *Invalid Routes Corruption* for route parameters by using non-existent airport codes or identical from/to locations;
- *Non-existent Booking Corruption* for functions like `cancel_booking` by generating random non-existent IDs.

**TradingBot** For the stock trading API, we implemented three corruption strategies:

Tool Name	Argument	Description	Domain Type	Domain Values	Data Dep.	Required
duplicate	output_filename	Name of duplicate file	string	Any filename	N	Y
rename	output_filename	New filename	string	Any filename	N	Y
search	object_name	Search term/object	string	Any search term	N	Y
count_pages	No arguments					
compress_file	output_filename	Compressed output name	string	Any filename	N	N
convert	format	Target format	finite	[pptx, doc, png, jpeg, tiff]	N	Y
	output_filename	Output filename	string	Any filename	N	Y
	zip	Zip output files	boolean	[true, false]	N	N
add_comment	page_num	Page number	numeric_range*	[1, num_pages]	Y	Y
	coordinates	Comment position [x,y]	list	[x, y] coordinates	N	Y
	font_size	Font size (points)	numeric_range	[8, 72]	N	Y
redact_page_range	start	Start page (inclusive)	numeric_range*	[1, num_pages]	Y	Y
	end	End page (inclusive)	numeric_range*	[1, num_pages]	Y	Y
redact_text	start	Start page	numeric_range*	[1, num_pages]	Y	Y
	end	End page	numeric_range*	[1, num_pages]	Y	Y
	object_name	Text to redact (list)	list	List of text strings	N	Y
	overwrite	Overwrite original	boolean	[true, false]	N	Y
	output_pathname	Output filename	string	Any filename	N	N
highlight_text	start	Start page	numeric_range*	[1, num_pages]	Y	Y
	end	End page	numeric_range*	[1, num_pages]	Y	Y
	object_name	Text to highlight (list)	list	List of text strings	N	Y
	overwrite	Overwrite original	boolean	[true, false]	N	Y
	output_pathname	Output filename	string	Any filename	N	N
underline_text	start	Start page	numeric_range*	[1, num_pages]	Y	Y
	end	End page	numeric_range*	[1, num_pages]	Y	Y
	object_name	Text to underline (list)	list	List of text strings	N	Y
	overwrite	Overwrite original	boolean	[true, false]	N	Y
	output_pathname	Output filename	string	Any filename	N	N
extract_pages	start	Start page	numeric_range*	[1, num_pages]	Y	Y
	end	End page	numeric_range*	[1, num_pages]	Y	Y
	overwrite	Overwrite original	boolean	[true, false]	N	Y
	output_pathname	Output filename	string	Any filename	N	N
delete_page	page_num	Page to delete	numeric_range*	[1, num_pages]	Y	Y
	overwrite	Overwrite original	boolean	[true, false]	N	Y
	output_pathname	Output filename	string	Any filename	N	N
delete_page_range	start	Start page	numeric_range*	[1, num_pages]	Y	Y
	end	End page	numeric_range*	[1, num_pages]	Y	Y
	overwrite	Overwrite original	boolean	[true, false]	N	Y
	output_pathname	Output filename	string	Any filename	N	N
add_signature	page_num	Page for signature	numeric_range*	[1, num_pages]	Y	Y
	position	Signature position	finite	[top-left, top-middle, ...]	N	Y
	overwrite	Overwrite original	boolean	[true, false]	N	Y
	output_pathname	Output filename	string	Any filename	N	N
add_page_with_text	text_content	Page text content	string	Any text content	N	Y
	font_size	Text font size	numeric_range	[8, 72]	N	Y
	page_num	Insert position	numeric_range*	[1, num_pages+1]	Y	Y
add_watermark	watermark_text	Watermark text	string	Any text	N	Y
	transparency	Transparency level	numeric_range	[0.0, 1.0]	N	Y
add_password	password	PDF password	string	Any password string	N	Y

Table 6: Document Plugin API: Complete Tool and Argument Specification with Domain Dependencies

Tool Name	Argument	Description	Domain Type	Domain Values	Data Dep.	Required
get_current_time				No arguments		
update_market_status	current_time_str	Time in HH:MM AM/PM	string	HH:MM AM/PM format	N	Y
get_symbol_by_name	name	Company name	string	Any company name	N	Y
get_stock_info	symbol	Stock symbol	string*	Available stock symbols	Y	Y
get_order_details	order_id	Order identifier	numeric_range*	Existing order IDs	Y	Y
cancel_order	order_id	Order to cancel	numeric_range*	Existing order IDs	Y	Y
place_order	order_type	Buy or Sell	finite	[Buy, Sell]	N	Y
	symbol	Stock symbol	string*	Available stocks	Y	Y
	price	Price per share	numeric_range	[0.01, 10000.0]	N	Y
	amount	Number of shares	numeric_range	[1, 10000]	N	Y
make_transaction	txct_type	Transaction type	finite	[deposit, withdrawal]	N	Y
	amount	Transaction amount	numeric_range	[0.01, 1000000.0]	N	Y
get_account_info				No arguments		
fund_account	amount	Funding amount	numeric_range	[0.01, 1000000.0]	N	Y
remove_stock_from_watchlist	symbol	Stock to remove	string*	Watchlist stocks	Y	Y
get_watchlist				No arguments		
get_order_history				No arguments		
get_transaction_history	start_date	Start date filter	string	YYYY-MM-DD format	N	N
	end_date	End date filter	string	YYYY-MM-DD format	N	N
update_stock_price	symbol	Stock symbol	string*	Available stocks	Y	Y
	new_price	New stock price	numeric_range	[0.01, 10000.0]	N	Y
get_available_stocks	sector	Market sector	finite	[Technology, Automobile, Healthcare, Finance, Energy]	N	Y
filter_stocks_by_price	stocks	Stock list to filter	list	List of stock symbols	N	Y
	min_price	Minimum price	numeric_range	[0.01, 10000.0]	N	Y
	max_price	Maximum price	numeric_range	[0.01, 10000.0]	N	Y
add_to_watchlist	stock	Stock to add	string*	Available stocks	Y	Y
notify_price_change	stocks	Stocks to monitor	list	List of stock symbols	N	Y
	threshold	Change threshold (%)	numeric_range	[0.01, 100.0]	N	Y

Table 7: Trading Plugin API: Complete Tool and Argument Specification with Domain Dependencies

Tool Name	Argument	Description	Domain Type	Domain Values	Data Dep.	Required
startEngine	ignitionMode	Engine ignition mode	finite	[START, STOP]	N	Y
fillFuelTank	fuelAmount	Fuel to add (gallons)	numeric_range*	[0, 50-current_fuel]	Y	Y
lockDoors	unlock	Lock or unlock	boolean	[true, false]	N	Y
	door	Doors to operate	list*	[driver, passenger, rear_left, rear_right]	Y	Y
adjustClimateControl	temperature	Target temperature	numeric_range	[-10, 50]	N	Y
	unit	Temperature unit	finite	[celsius, fahrenheit]	N	N
	fanSpeed	Fan speed (0-100)	numeric_range	[0, 100]	N	N
	mode	Climate mode	finite	[auto, cool, heat, defrost]	N	N
get_outside_temperature_from_google				No arguments		
get_outside_temperature_from_weather_com				No arguments		
setHeadlights	mode	Headlight mode	finite	[on, off, auto]	N	Y
displayCarStatus	option	Status display option	finite	[fuel, battery, doors, climate, headlights, parkingBrake, brakePedal, engine]	N	Y
activateParkingBrake	mode	Brake mode	finite	[engage, release]	N	Y
pressBrakePedal	pedalPosition	Pedal position (0-1)	numeric_range	[0, 1]	N	Y
releaseBrakePedal				No arguments		
setCruiseControl	speed	Cruise speed (mph)	finite*	[0, 5, 10, ..., 120]	Y	Y
	activate	Activate cruise	boolean*	[true, false]	Y	Y
	distanceToNextVehicle	Following distance (m)	numeric_range	[0, 1000]	N	Y
get_current_speed				No arguments		
display_log	messages	Log messages	list	List of strings	N	Y
estimate_drive_feasibility_by_mileage	distance	Distance in miles	numeric_range	[0, 10000]	N	Y
liter_to_gallon	liter	Liters to convert	numeric_range	[0, 1000]	N	Y
gallon_to_liter	gallon	Gallons to convert	numeric_range	[0, 1000]	N	Y
estimate_distance	cityA	First city zipcode	finite	[83214, 74532, 56108, ...]	N	Y
	cityB	Second city zipcode	finite	[83214, 74532, 56108, ...]	N	Y
get_zipcode_based_on_city	city	City name	finite	[Rivermist, Stonebrook, ...]	N	Y
set_navigation	destination	Destination address	string	Street, city, state format	N	Y
check_tire_pressure				No arguments		
find_nearest_tire_shop				No arguments		

Table 8: Vehicle Control Plugin API: Complete Tool and Argument Specification with Domain Dependencies

Tool Name	Argument	Description	Domain Type	Domain Values	Data Dep.	Required
pwd	No arguments					
ls	a	Show hidden files	boolean	[true, false]	N	N
cd	folder	Directory to change to	string*	Available directories + [..., /]	Y	Y
mkdir	dir_name	New directory name	string	Any valid directory name	N	Y
touch	file_name	New file name	string	Any valid filename	N	Y
echo	content	Text content	string	Any text string	N	Y
	file_name	Output file (optional)	string	Any filename	N	N
cat	file_name	File to display	string*	Available files	Y	Y
find	path	Search starting point	string	Any path	N	N
	name	Search pattern	string	Any search pattern	N	N
wc	file_name	File to count	string*	Available files	Y	Y
	mode	Count mode	finite	[l, w, c]	N	N
sort	file_name	File to sort	string*	Available files	Y	Y
grep	file_name	File to search	string*	Available files	Y	Y
	pattern	Search pattern	string	Any text pattern	N	Y
du	human_readable	Human readable format	boolean	[true, false]	N	N
tail	file_name	File to display	string*	Available files	Y	Y
	lines	Number of lines	numeric_range	[1, 100]	N	N
diff	file_name1	First file	string*	Available files	Y	Y
	file_name2	Second file	string*	Available files	Y	Y
mv	source	Source file/directory	string*	Available items	Y	Y
	destination	Destination name	string*	Available items + new names	Y	Y
rm	file_name	File/directory to remove	string*	Available items	Y	Y
rmdir	dir_name	Directory to remove	string*	Available directories	Y	Y
cp	source	Source file/directory	string*	Available items	Y	Y
	destination	Destination name	string*	Available items + new names	Y	Y

Table 9: File System Plugin API: Complete Tool and Argument Specification with Domain Dependencies

- *Invalid Symbols Corruption* for functions like `get_stock_info` by using non-existent symbols or malformed formats;
- *Financial Validation Corruption* for `place_order` and related functions by using negative values or amounts exceeding account balance;
- *Order State Conflicts Corruption* for `cancel_order` operations by referencing completed orders or using malformed order IDs.

Plugin	Update Trigger	Dynamic Domain Updates	Affected Operations
<b>Travel</b>			
	Credit card registration	Card IDs → available payment methods	book_flight, get_credit_card_balance, purchase_insurance
	Flight booking	Booking IDs → cancellable/retrievable bookings	cancel_booking, retrieve_invoice, re-contact_customer_support
	Budget setting	Budget limits → financial constraints	All cost-related operations
	Route updates	Airport codes → valid travel routes	get_flight_cost, book_flight
<b>Document</b>			
	Page operations	Page count → valid page numbers	All page-specific operations
	Document loading	Total pages → range constraints	add_comment, delete_page, etc.
	Cache invalidation	State changes → domain refresh	Page-changing operations
<b>Trading</b>			
	Order placement	Order IDs → manageable orders	get_order_details, cancel_order
	Stock updates	Available stocks → tradeable symbols	place_order, get_stock_info
	Watchlist changes	Watchlist → removable stocks	remove_stock_from_watchlist
<b>Vehicle</b>			
	Fuel level changes	Current fuel → addable amount	fillFuelTank
	Door state changes	Door status → operable doors	lockDoors
	Engine state	Running/stopped → cruise control availability	setCruiseControl
<b>File System</b>			
	Directory navigation	Current contents → available items	cd, cat, mv, cp, rm
	File operations	File list → operable files	File-specific operations
	Directory changes	Directory list → navigable paths	cd, rmdir
	State synchronization	FS changes → domain cache invalidation	All state-changing operations

Table 10: Dynamic Domain Update Rules and Triggers Across Plugin System

## Human Annotation Guidelines

### Objective:

Annotators must evaluate five LLM-generated queries per sample. Each query is scored on three dimensions: (A) Naturalness of language, (B) Faithfulness to the expected tool call, and (C) Executability/Validity. Additionally, annotators must check for removal of Personally Identifiable Information (PII), assess tool call feasibility, and select one optimal query per sample.

### Evaluation Rubric

Criterion	Score 5	Score 4	Score 3	Score 2	Score 1
<b>A. Naturalness</b>	Fully fluent, natural, human-like	Minor awkwardness or stiffness	Understandable but robotic	Clearly awkward or difficult to read	Unintelligible or nonsensical
<b>B. Faithfulness</b>	Perfect match to expected tool call; all required arguments present	Mostly aligned; minor phrasing or parameter issues	Some omissions or hallucinations; core logic intact	Major deviations from expected tool behavior	Entirely incorrect or misleading tool structure
<b>C. Executability</b>	Fully executable; properly structured and valid	Executes with minor issues or missing defaults	Partially executable with moderate corrections needed	Major issues preventing execution	Unexecutable or contradicts tool logic/API

### Required Checks

- **PII Removal:** Ensure no personal identifiers (names, emails, phone numbers, IDs) are present. Flag these queries for further processing.
- **Tool Call Validation:** If feasible, simulate or run tool calls to confirm validity and argument correctness.
- **Error Identification:** Mark and annotate any queries with logical inconsistencies, invalid parameters, or unsupported constraints.

Figure 7: Summary of instructions given to human annotators.