au^2 -Bench: Evaluating Conversational Agents in a Dual-Control Environment

Anonymous authors

000

001

003 004

010 011

012

013

014

015

016

018

019

021

024

025

026

027

028

031

032

034

037

038

040 041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Existing benchmarks for conversational AI agents simulate single-control environments, where only the AI agent can use tools to interact with the world, while the user remains a passive information provider. This differs from real-world scenarios like technical support, where users need to actively participate in modifying the state of the (shared) world. In order to address this gap, we introduce τ^2 -bench, with four key contributions: 1) A novel Telecom dual-control domain modeled as a Dec-POMDP, where both agent and user make use of tools to act in a shared, dynamic environment that tests both agent coordination and communication, 2) A compositional task generator that programmatically creates diverse, verifiable tasks from atomic components, ensuring domain coverage and controlled complexity, 3) A reliable user simulator tightly coupled with the environment, whose behavior is constrained by tools and observable states, improving simulation fidelity, 4) fine-grained analysis of agent performance through multiple ablations including separating errors arising from reasoning vs communication/coordination. In particular, our experiments show significant performance drops when agents shift from no-user to dual-control, highlighting the challenges of guiding users. Overall, τ^2 -bench provides a controlled testbed for agents that must both reason effectively and guide user actions.

1 Introduction

Existing benchmarks for conversational AI agents are designed to test their abilities to communicate effectively with a user and perform the right sequence of actions to solve tasks (Yao et al., 2024; Lu et al., 2024; Xiao et al., 2024; Prabhakar et al., 2025). These benchmarks are inherently single-control environments, where the AI agent is able to interact with the world but the (simulated) user is limited to providing information about preferences and goals. In τ -bench (Yao et al., 2024) for example, the retail and airline domains test the agent's ability to solve constraint satisfaction tasks, where constraints stem from a combination of a domain policy that the agent must follow and the user's cognitive state including beliefs, goals and preferences.

In such settings, the user's understanding of the agent's environment (including the set of actions the agent can take) is provided through carefully crafted natural language instructions that help ensure a single, solvable path for the constraint satisfaction problem. While this enables easier task specification, the user can only perceive the agent's actions through communication and reason about the environment state solely based on the initial instructions. This is quite different from real-world scenarios like technical customer support, where the user has to actively participate in taking some actions to diagnose and solve problems when asked to do so, such as *restarting their phone* or *turning off airplane mode*.

In order to capture this additional real-world complexity, we introduce τ^2 -bench with a *dual-control* environment, where (LLM-simulated) users can take actions and call tools in addition to communicating with the agent. The dual-control model provides a few advantages in improving the user simulation, such as selective information hiding, non-verbal manipulation of the environment, and easier specification of the user's task or scenario, especially in cases where the structured format of tools is preferred to unstructured natural language.

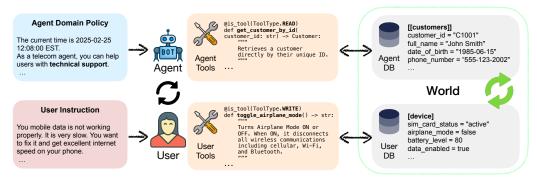


Figure 1: Supporting dual-control environment in τ^2 -bench. The agent has access to a set of tools that interact with a database, and is tasked with resolving the user's request via Tool-Agent-User (TAU) interactions while adhering to the domain policy. To test real-world scenarios, the user is simulated by another AI agent given a scenario-based instruction and a set of tools that interact with its own database. The simulated user can be regarded as handling an easier version of the TAU interaction in a dual format (Tool-User-Agent), where it only need to follow instructions but does not need to reason about solutions for the task.

The primary challenge in a dual-control setup lies in granting the user simulator meaningful agency via tools while maintaining the "complexity asymmetry" between the agent and the user simulator – that is, making sure the user's affordances are limited and the user still requires support from the agent in order to solve issues. While tools for the user simulator could disrupt this balance by increasing its capabilities, we find that carefully designing the tools can help constrain user behavior, leading to more reliable simulations. We achieve this by ensuring user tools yield only human-readable outputs, limiting user planning to reactive tool use based on agent requests, and tightly constraining user behavior through the environment. This results in a controlled testbed for agents that must both think and effectively guide a user to take the right actions.

Overall, τ^2 -bench makes four key contributions:

Telecom dual-control domain. We introduce a novel dual-control telecom domain within τ^2 -bench, where, unlike previous benchmarks, both the AI agent and the user possess distinct tools to observe, act upon, and verify the state of a shared, dynamic environment. This is formalized using a Decentralized Partially Observable Markov Decision Process (Dec-POMDP) (Oliehoek et al., 2016). This dual-control setup is designed to accurately represent real-world collaborative scenarios and exposes crucial agent coordination and communication challenges absent from evaluations where users have limited agency. Experiments show state-of-the-art LLMs struggle significantly in this domain (e.g., pass^1 of 34% for gpt-4.1 (OpenAI, 2025a), 42% for o4-mini (OpenAI, 2025b), and 49% for claude-3.7-sonnet (Anthropic, 2025) on new tasks).

Compositional task generator. τ^2 -bench incorporates a programmatic task generator that automatically composes a vast and diverse set of verifiable tasks from a small set of atomic base scenarios (defined by initialization, solution, and assertion functions). This method ensures provable correctness of tasks, provides complete domain coverage, allows for explicit control over task complexity (e.g., by number of solution steps or issue type), and removes the manual effort and potential brittleness associated with hand-crafted task suites.

Reliable user simulator. We enhance the reliability of the user simulation by tightly coupling the user simulator to the environment. User behavior is constrained by the available tools and the observable state of the environment, leading to more predictable and consistent interactions. This approach significantly alleviates the need for complex natural language prompting to guide the user simulator and results in substantially higher reliability (e.g., the telecom domain's user simulator shows a 16% error rate with 6% critical errors, compared to 40% error rate with 12% critical errors in the retail domain from τ -bench).

Fine-grained diagnosis of agent failures. τ^2 -bench enables a decomposed diagnosis of agent performance by evaluating task success in different modes: (i) a fully autonomous mode ("no-user mode"), where the agent controls all tools, isolating its reasoning capabilities, and (ii) the standard dual-control mode, which introduces communication and coordination requirements. Our findings reveal a substantial performance decrease (around 20% pass 1) when agents must shift from autonomous operation to guiding a user. This clearly distinguishes pure reasoning failures from those arising from communication and decentralized control, pinpointing the latter as a critical bottleneck.

2 RELATED WORK

Benchmarks for Conversational AI Agents. Following a long line of research into language agents (Yao et al., 2022; Zhou et al., 2023; Jimenez et al., 2023; Liu et al., 2023; Ruan et al., 2023), LLM tool use (Yan et al., 2024; Qin et al., 2024; Huang et al., 2023), and task-oriented dialog (Chen et al., 2021; Budzianowski et al., 2018; Andreas et al., 2020; Schatzmann et al., 2007; Gür et al., 2018; He et al., 2018; Hu et al., 2023), τ -bench (Yao et al., 2024) is a recently introduced benchmark to measure the reliability of language agents in multi-turn task-oriented conversations such as customer service workflows, while respecting domain rules. Each task in τ -bench instantiates a live conversation between a user simulator and the language agent, with tasks spread across two domains – retail and airline. To quantify reliability the paper introduces pass k metrics: the fraction of k independent runs that succeed.

Several follow-ups to τ -bench have explored variations of the basic setup. FlowBench (Xiao et al., 2024) isolates the planning step of tool-using agents by injecting explicit workflow knowledge into the prompt, using natural language, python-like pseudocode or mermaid flowcharts. IntellAgent (Levi & Kadar, 2025) provides an evaluation pipeline to programmatically build synthetic test suites from structured policy graphs that encode domain rules and their co-occurrence statistics. IntellAgent explicitly uses τ -bench as an external gold standard, reporting a high Spearman correlation between the two score distributions, and acts as a fast, synthetic proxy task. APIGen-MT (Prabhakar et al., 2025) explores the idea of fine-tuning tool-calling agents for τ -bench. They generate data by creating conversation blueprints which are sequences of tool calls that depend on each other, followed by simulating conversational traces based on each blueprint. ToolSandbox (Lu et al., 2024) focuses on creating stateful tools in order to evaluate agent progress in a more fine-grained manner.

Our work extends the τ -bench paradigm and generalizes it to allow for both the user and agent to have state-changing abilities (via tool calls) over a shared world. As demonstrated in the results, this allows us to build more complex domains to test conversational agents, while also providing the opportunity for fine-grained analysis of agent failure points that can be improved upon.

User Simulation for Conversational Agents. The reliability of user simulation has been a key concern for benchmarks like τ -bench (Yao et al., 2024). While most efforts have focused on introducing supervision for the user simulator, for instance, by using a generic LLM to generate or validate user responses (Prabhakar et al., 2025), less attention has been paid to the possibility of using the environment itself to constrain and shape user simulator behavior for increased reliability, a core tenet of our approach. This concern has been extensively studied in the context of task-oriented dialogue systems, with early work by (Pietquin & Hastie, 2013) providing a comprehensive survey of metrics for evaluating user simulations. More recently, (Kazi et al., 2024) has demonstrated how LLMs can be effectively used as user-agents for evaluating task-oriented dialogue systems, showing that careful prompting and state tracking can lead to more reliable and context-aware user simulations.

Multi-Agent Benchmarks. Our work is also related to efforts to build multi-agent frameworks and evaluate them (Zhu et al., 2025). While we can consider the user and the agent in our paper as forming a multi-agent system, the key difference in our case is that the final evaluation still focuses on the agent's ability to elicit the right information from the user and perform the correct actions to solve the task. This introduces an inherent asymmetry between the agent and the user; our focus is not on solving a pure multi-agent problem but rather on the agent's capability to effectively guide and collaborate with a user who also possesses agency. In this sense, the framework can be collaborative (e.g., troubleshooting), competitive (e.g., negotiating a subscription), or a hybrid, requiring the agent to identify and navigate the scenario appropriately, even accounting for user mistakes or errors.

au^2 -венсн: Evaluating Agents in a Dual-Control Environment

 au^2 -bench serves as a platform for a systematic study of multi-turn interactions between a conversational AI agent and a simulated user. Dual-control interactions are modeled as a Decentralized Partially Observable Markov Decision Process (Dec-POMDP) (Oliehoek et al., 2016) where both the agent and user can communicate, use tools, and receive observations. This allows us to simulate complex scenarios like technical troubleshooting where agent and user must coordinate their actions to solve the task.

Table 1: Key statistics for the τ^2 -bench domains.

	retail	airline	telecom
Agent Databases	500 users, 50 products, 1,000 orders	500 users, 300 flights, 2,000 reservations	5 plans, 9 lines, 4 customers
Agent Tools	7 write, 6 read	6 write, 6 read	6 write, 7 read
User Tools	-	-	15 write, 15 read
Tasks	115	50	114 (full: 2285)

3.1 THE DEC-POMDP FORMALISM

As illustrated in Figure 1, the Dec-POMDP in τ^2 -bench involves two players: an agent and a user. The entire process is formally defined by the tuple $(S, \{A_i\}, \{\mathcal{O}_i\}, \mathcal{T}, \mathcal{R}, \mathcal{U}, \mathcal{M})$, where $i \in \{agent, user\}$ denotes players and each component in the tuple is detailed below, with illustrative examples drawn from the new telecom domain.

Message space (\mathcal{M}): The set of all possible (natural language) messages exchanged between the agent and the user. For example, the user could say "I cannot use mobile data." and the agent could respond with "Could you check whether your airplane mode is on?"

State space (S): The global state $S = S_{world} \otimes S_{history}$, where $S_{world} = S_{db,agent} \otimes S_{db,user}$ represents the *underlying* database states for the agent and user, and $S_{history}$ logs all interaction events (actions, observations, messages). For example, in the telecom domain, $S_{db,agent}$ might be CRM data (customer profiles, lines), while $S_{db,user}$ could be phone status.

Action spaces (\mathcal{A}_i): Player i's action $a_i \in \mathcal{A}_i$ is either a tool call $a_{i,tool} \in \mathcal{A}_{i,tool}$ (interacting with $\mathcal{S}_{db,i}$ via function calls like tool_name(**kwargs)) or a message $m_i \in \mathcal{M}$. Only one player acts per turn. In the telecom domain, the agent can access tools like get_customer_by_id(id) and the user can access tools like toggle_airplane_mode().

Observation spaces (\mathcal{O}_i): Player i's observation $o_i \in \mathcal{O}_i$ is either a tool observation $o_{i,tool}$ (e.g., data, messages, or errors from $a_{i,tool}$) or a message $m_j \in \mathcal{M}$ from player $j \neq i$. Only one player receives an observation per turn. In the telecom domain, the agent might observe customer details from get_customer_by_id, and the user might observe a message indicating the airplane mode has been turned off from toggle_airplane_mode.

Transition function (\mathcal{T}): Defines system dynamics via $\mathcal{T}: \mathcal{S} \times \mathcal{A} \to \mathcal{S} \times \mathcal{O}$. Given current state $s \in \mathcal{S}$ and joint action $a = (a_{agent}, a_{user})$, it yields a new state $s' \in \mathcal{S}$ and joint observation $o = (o_{agent}, o_{user})$. Calling a tool $a_{i,tool} \in \mathcal{A}_{i,tool}$ may change \mathcal{S}_{world} and yield $o_i \in \mathcal{O}_{i,tool}$. Sending a message $m_i \in \mathcal{M}$ yields $o_j = m_i$ for $j \neq i$. In both cases, s' includes updated \mathcal{S}_{world} and $\mathcal{S}_{history}$. For example, an agent's action enable_roaming (customer_id, line_id) would update the world state (the roaming service for the specific line number is enabled), and a user's action toggle_airplane_mode would update the status of the mocked phone.

Reward function (\mathcal{R}): A function $\mathcal{R}: \mathcal{S} \to [0,1]$ providing a global reward based on the overall state $s \in \mathcal{S}$ (database states, history), signaling task success or failure. For example, in telecom, the agent is rewarded if the issue ("no mobile data") is fixed, as verified by the user's database state.

Instruction space (\mathcal{U}): The instruction space \mathcal{U} defines the scenario guiding realistic user simulation, as well as the domain policies to which the agent must adhere when assisting the user.

The Dec-POMDP formalism offers key advantages for simulating complex, interactive scenarios (see Figure 2 for an example trajectory of interactions). It enables realistic simulations of collaborative environments, such as technical support, where users perform actions guided by agents. This presents agents with crucial coordination and communication challenges. In addition, the formalism enhances the reliability and control of user simulation. By predefining user tools and their effects on user states, user behavior becomes more controllable and less reliant on extensive natural language prompting.

3.2 Domain and task creation

Similar to τ -bench, we adopted a multi-stage creation process to build domain-specific materials for new domains. This process, illustrated using the telecom domain, involves the following stages:

Stage 1: Creating agent's database schema and tools. We begin by prompting Large Language Models (LLMs) to generate a Product Requirements Document (PRD) that outlines the domain's core business logic. This PRD specifies the database schema and necessary functions. In the telecom

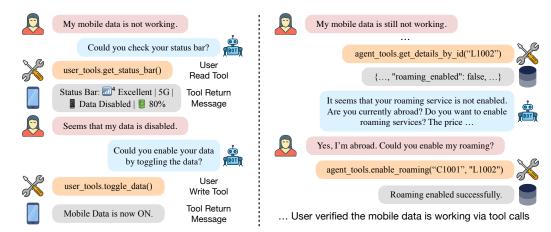


Figure 2: An example agent-user interaction trajectory ($S_{history}$) of τ^2 -bench in the telecom domain. By controlling the implementation of the user tools (the mocked phone), we can reliably simulate the user's response to agent's actionable instructions like "checking the status bar" and "toggling data" based on the underlying database state. On the right half, we show the possibility of modeling the impact of agent's tool calls on the user's database state, where the roaming service for the user is enabled on the agent's end and therefore allows the user's phone to roam.

domain, this involved defining a customer CRM system with schemas for customers and lines, along with functions to manage them. An LLM then generates function implementations, a mock database, and unit tests based on the PRD. We manually refine the generated code until all unit tests pass.

Stage 2: Creating user's database schema and tools. For troubleshooting scenarios, we similarly use an LLM to define the user's database schema and tools. In the telecom domain, this included implementing a mocked user phone device with status (e.g., signal strength) and functions (e.g., toggling airplane mode). Again, an LLM generates implementations, a mock database, and unit tests, which are then manually refined until all tests pass.

Stage 3: Programmatic task creation. We employ a combinatorial approach to generate diverse, verifiable tasks from atomic building blocks (see Appendix A.3 for details on our task factorization framework).

Each atomic subtask t is about a specific problem to be resolved, for example, airplane mode on leads to mobile data not working. Specifically, each subtask t is defined as $(\{f_{t,k}^{init}\}, \{f_{t,k}^{sol}\}, \{f_{t,k}^{assert}\})$, where $f_{t,k}$ is the k-th function call of the subtask t that interacts with the agent's or user's database:

- Initialization functions $f_{t,k}^{init}$ specify calls to set up the initial task state, typically by updating the database values. For instance, in telecom, an initialization might be set_airplane_mode(True).
- Solution functions $f_{t,k}^{sol}$ specify tool calls to resolve issues introduced by initialization. For example, toggle_airplane_mode() could be a solution for the initialization example given above. Note that these must be tools available to the agent or user.
- Assertion functions $f_{t,k}^{assert}$ specify conditions the final state $\mathcal S$ must meet for the task to be considered solved. For instance, assert_service_status("connected") checks if the user's service is active in telecom.

While solution functions $f_{t,k}^{sol}$ are restricted to agent or user tools, initialization and assertion functions can be any function in the relevant database.

Atomic subtasks are grouped such that mutually exclusive or alternative subtasks are in the same group. A composite task is created by selecting at most one subtask from each group, concatenating their respective function calls. Task correctness is automatically verified by checking if the final state $s \in \mathcal{S}$ satisfies all assertion functions after applying initialization and then solution functions. We also verify that the task is not resolved until all solution functions are applied.

In the telecom domain, we developed 15 atomic subtask groups for 3 user intents of increasing complexity: service_issue, mobile_data_issue, and mms_issue. Combining these subtasks in a programatic way yields 2285 tasks. We then subsample 114 tasks to form a balanced

distribution over different intents and numbers of subtasks (details in Appendix A). The number of subtasks in a task serves as a proxy for difficulty as more diagnostic and resolution steps are required.

Stage 4: Creating domain-specific agent policy. Based on the curated tasks and their solutions, we prompt LLMs to generate domain-specific policies for the agent. For troubleshooting, these policies guide the agent in diagnosing and resolving user issues, often outlining step-by-step procedures for common problems related to each user intent, details in Appendix D.2.

Stage 5: Manual refinement. We jointly refine all the domain materials including tools, policy and atomic subtasks to improve the quality of the domain.

Moreover, compared to the original τ -bench, τ^2 -bench enables a developer to specifically associate a task with a Persona (a brief description of the user's identity). We put this to use in our new domain. Each telecom domain task was randomly assigned one of the following personas: None, Easy, and Hard. The None persona means that no specific persona is provided to the user simulator. The Easy persona describe the profile of a user who is rather familiar with the domain while the Hard persona represents a more challenging user with low technical knowledge (see Appendix A.1).

3.3 TASK EVALUATION

The success of a task can be defined by different criteria: DB check, status assertions, natural language assertions, communication info check, and action matching. The DB check and communication info check are the same as the original τ -bench. The status assertion involves verifying specific conditions in the final world state \mathcal{S}_{world} using the predefined assertion functions (e.g., checking if a service is connected). The natural language assertion involves verifying specific conditions in the final history state $\mathcal{S}_{history}$ using a natural language description, like "the agent diagnosed the cause of the issue." The action matching involves verifying if every solution function $f_{t,k}^{sol}$ exists in the actual agent-user interaction trajectory. Practically, each task can specify a subset of these criteria based on its features. In telecom, only assertion functions are used to evaluate task success.

4 EXPERIMENTS

4.1 AGENT SETTINGS

All LLM API calls are implemented using the Litellm package (BerriAI, 2025). We evaluated four large language models: gpt-4.1-mini-2025-04-14, gpt-4.1-2025-04-14, o4-mini-2025-04-16, and claude-3-7-sonnet-20250219. The user simulator is implemented using gpt-4.1-2025-04-14. Each task is run four times, maintaining a consistent LLM temperature of 0 to promote deterministic outputs. Both the agents and the user simulator are implemented as function-calling agents. All tools are provided to LLMs in the OpenAI tools format. The agent prompt includes generic guidelines along with domain-specific policies. Similarly, the user prompt contains generic guidelines supplemented by task-specific instructions. Both domain policies and prompts are available in Appendices C and D.

When the gpt-4.1-2025-04-14 agent is paired with the gpt-4.1-2025-04-14 user simulator, the average agent/user simulation costs are 0.086/0.059 per task, respectively. The cost of running all domains for 1 trial per task is approximately \$40.

4.2 RESULTS

Pass^k scores. We computed performance metrics on the verified τ^2 -bench domains (retail and airline) and on our new telecom domain (see Figure 3). Our findings indicate that the telecom domain presents a greater challenge, exhibiting an overall lower success rate compared to other domains. gpt-4.1 pass^1 drops from 74%/56% for retail and airline respectively to 34% for telecom. gpt-4.1-mini, o4-mini, and claude-3.7-sonnet perform better with pass^1 of around 50% for telecom. In the case of claude-3.7-sonnet, the pass^1 score for telecom (49%) is on par with airline. However, as k increases, the pass^k scores decline more rapidly for telecom compared to airline, suggesting less consistent performance on the telecom domain.

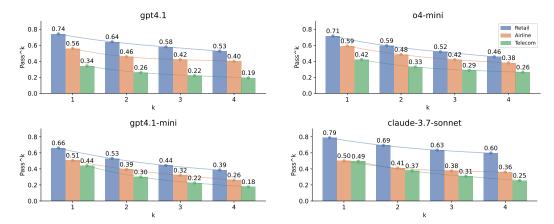


Figure 3: pass k metrics across all evaluated domains (airline, retail, telecom). **Top Left**: gpt-4.1, **Top Right**: o4-mini, **Bottom Left**: gpt-4.1-mini, **Bottom Right**: claude-3.7-sonnet.

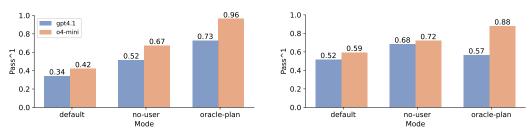


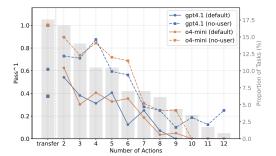
Figure 4: pass¹ metrics for the Telecom domain across different operational modes (Default, No-User, Oracle Plan) using the Default user simulation. **Left**: original policy. **Right**: workflow-based policy. This figure illustrates the impact of reasoning load and decentralized control on agent performance.

Ablation analysis. An agent's success in τ^2 -bench depends on two things -1) how well it can communicate and collaborate with the user at solving the issue, and 2) how well it can reason over and apply the domain guidelines specified in the policy document. In order to understand the impact of each of these components—reasoning and communication for dual control—we perform an ablation study. Specifically, we evaluate performance in the telecom domain across three distinct settings:

- **Default**: The default agent and user simulator configuration where the agent and user collaborate in a *dual-control* setup.
- **No-User**: The agent is provided with a ticket summarizing the user's problem and success criteria. The agent controls all tools, including those typically operated by the user, and is solely responsible for solving the problem. This setting tests the agent's reasoning and tool-calling capability independently of its capacity to interact with the user.
- Oracle Plan: The agent is provided with the sequence of tool calls required to solve the problem, encompassing actions for both the agent and the user. This setup alleviates the agent's reasoning load, focusing on its ability to collaborate with the user to execute a known plan.

Figure 4 (left) reports performance across these settings for gpt-4.1 and o4-mini, revealing key insights. The difference between Oracle Plan and Default configurations highlights the impact of the reasoning load on agent performance. Unsurprisingly, providing the ground truth leads to better performance than the Default setting. But it is notable that this effect is larger for the o4-mini than for gpt-4.1, suggesting that o4-mini is better able to make use of the ground truth information.

The comparison between No-User and Default modes illustrates the impact of dual control and the associated communication overhead on agent error rates and overall success. For both models, shifting from no user operation (No-User) to a collaborative setup (Default) where the agent must guide the user results in a substantial drop in pass^1 (18% drop for gpt-4.1 and 25% drop for o4-mini). This underscores that LLMs still face significant challenges when solving problems with an active user who shares control of the environment.



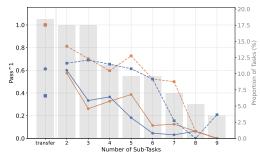
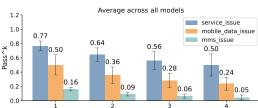


Figure 5: pass¹ scores across various tasks in telecom binned by the number of actions required to solve them (left) or the number of different issues that need to be addressed (right). *transfer* refers to the special case of a task that requires to be transferred to a human and cannot be solved by the agent alone. (Grey bars indicates the proportion of the tasks that fall into that bin.)



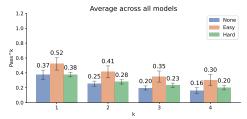


Figure 6: pass^k per issue (left) and persona (right) type for the telecom domain aveaged across all models. Performance is shown for service_issue, mobile_data_issue, and mms_issue issue types and across different persona types, highlighting how different issue and persona types affect success rates.

Impact of policy document on performance. Figure 4 also allows us to analyze the impact of the policy document on task success. Specifically, we created an alternate policy document that provides more specific details on the workflow required to solve each type of task, with the rationale that having the workflow provides more concrete guidance to the agent (see Appendices D.2.2 and D.2.3). We observe that this is indeed the case and slightly improves agent performance using the workflow policy (right) over the original one (left) under the Default and No-User modes. Surprisingly, workflow policy document hurts performance on Oracle Plan for both gpt-4.1 and o4-mini. Our hypothesis here is that since the agent already has the ground truth action sequence, providing it the workflow might lead to confusion and hurt its performance rather than help.

Impact of number of actions and sub-tasks. Figure 5 breaks down the pass¹ scores across various tasks in telecom binned by the number of actions required to solve them (left) or the number of different sub-tasks that need to be addressed (right). *transfer* refers to the special case of a task that requires to be transferred to a human and cannot be solved by the agent alone.

As expected, regardless of the base model being $\mathtt{gpt-4.1}$ or $\mathtt{o4-mini}$, agent performance drops as the number of actions increases, reaching close to zero for > 7 actions in <code>Default</code> mode. The <code>No-User</code> mode results in higher scores overall, although the gap reduces (from about 0.3-0.4 to < 0.2) as the number of actions increases. This hints that maintaining reliability over longer-horizon tasks remains a challenge under both settings and communication with the user is not the only bottleneck. Interestingly, for <code>No-User</code>, <code>gpt-4.1</code> performs better at the tail end (10 actions or more) than o4-mini.

We observe a similar trend with increase in the number of distinct sub-tasks per task – performance trends downwards for both base models, with the no-user mode generally being higher than the default mode. Both these results validate that our domain design and task creation process provide a natural path to scaling complexity via combining different sub-tasks into a single task.

Impact of issue and persona types. The telecom domain is organized around three primary user intents and personas reflecting different difficulties. Figure 6 provides a breakdown analysis of performance by issue and persona types. We observe that the pass^k scores for all LLMs are driven by higher failure rates on more complex issue types (mobile_data_issue and mms_issue), while the agent tends to perform better on tasks associated with the Easy persona compared to those associated with the Hard and (interestingly) the None one. Refer to Appendix A.2 for more details.

Table 2: User simulator error count (*rates*) across domains. Only critical user errors prevent the agent from solving the task while benign errors do not affect task completion. New telecom domain shows much lower error rate and no critical errors were reported. (See Appendix E for more details).

Domain	Num Conversations	Critical Errors	Benign Errors	Total Errors
airline	100	13 (13%) 6 (12%) 3 (6%)	34 (34%)	47 (47%)
retail	50	6 (12%)	14 (28%)	20 (40%)
telecom	50	3 (6%)	5 (10%)	8 (16%)

4.3 How does dual-control impact benchmark reliability?

Ensuring the reliability of conversational agent benchmarks is paramount. Three primary sources of uncertainty can impact benchmark reliability: **implementation errors**, **task specification errors**, and **user simulator errors**. While the user simulator is often cited as a critical component requiring careful evaluation, its assessment can be confounded by issues in the benchmark's implementation or task definitions. Therefore, we first address these potential error sources before evaluating the user simulator itself. This is detailed in Appendix B.

User simulator quality evaluation. Having minimized errors in the benchmark implementation and task specifications, we evaluated the user simulator quality. To assess the quality of the user simulator across domains, we manually annotated interaction traces generated using gpt-4.1 for both the User Simulator and the Agent. Each conversation was reviewed by two separate annotators tasked with identifying user simulator errors. Annotators were given the User Simulator Guidelines, the specific User Instructions for this conversation (see Appendix C.2), descriptions of the available User Tools (if any), and the complete conversation trajectory (messages and tool calls). Annotators assessed each user turn against four criteria: adherence to User Simulator Guidelines, adherence to User Instructions, correct use of User Tools, and generation of a natural and consistent conversational continuation. Errors were categorized as either (1) task-critical errors: high-severity failures that preclude task completion (e.g., generating an intent that contradicts the user goal, or causing an irrecoverable state transition), or as (2) task-benign errors: Errors that do not prevent the task from being completed.

Reliability of the user simulator. As shown in Table 2, our analysis of user simulator behavior reveals significant improvements in reliability for the new telecom domain. While for the retail and airline domains we recorded a 40% and 47% error rate for the user simulator (with 12% and 13% being critical errors that prevent task completion), this rate is much lower for the telecom domain, only 16% with 6% critical errors reported. This substantial improvement in reliability can be attributed to the domain design, which shapes and tightly constrains user behavior through its environment and available affordances. Rather than relying heavily on natural language specifications to guide behavior, the telecom domain's structured interface and clear action space naturally guide the user simulator toward correct interactions, resulting in more consistent and predictable behavior.

5 CONCLUSION

We present τ^2 -bench that generalizes τ -bench by introducing the dual-control setting and found a substantial performance drop in LLMs due to coordination and communication requirements, highlighting these as critical bottlenecks over pure reasoning capabilities for solving user requests.

More work remains to be done to improve the user simulator. Although we have shown that augmenting users with curated tools can help avoid critical errors, we have not yet investigated how this method could be applied to the existing airline and retail domains. Doing so would pave the way towards a more generic solution to ensuring high quality user simulator. Extending domain coverage for the benchmark still heavily relies on human experts. For benchmarking methods to be adopted by industry, providing much needed standards, it is critical to further investigate how to automate the domain curation process.

One important limitation of τ^2 -bench is that it does not explicitly model the expert-novice gap inherent to most customer support tasks. When interacting with a naive user, an expert must understand the user's mental model and adapt explanations accordingly. Assessing and improving the AI agent's abilities to bridge this gap is a promising direction for future work and τ^2 -bench provides a strong starting point for such explorations.

REPRODUCIBILITY STATEMENT

We have provided all codes and data in the supplementary materials and will make them open source. The experiments can be reproduced by following the experimental settings and our README in the code.

ETHICS STATEMENT

The development of standardized benchmarks for Large Language Models (LLMs) and AI agents is crucial for ensuring societal control and fostering fairness amidst rapidly advancing technologies. Such benchmarks not only provide a framework for transparent evaluation but also enable research groups to coordinate their efforts around common tasks, thereby accelerating the overall progress in the field. While this work itself may not have direct negative societal implications, it contributes to the development of real-world agents, which will invariably have diverse economic and societal consequences. Therefore, it is also of paramount importance that AI agents are designed to collaborate effectively and safely with human users, a prerequisite for their responsible integration into commercial settings and everyday life.

REFERENCES

- Jacob Andreas, John Bufe, David Burkett, Charles Chen, Josh Clausman, Jean Crawford, Kate Crim, Jordan DeLoach, Leah Dorner, Jason Eisner, et al. Task-oriented dialogue as dataflow synthesis. *Transactions of the Association for Computational Linguistics*, 8:556–571, 2020.
- Anthropic. Claude 3.7 Sonnet, 2025. URL https://www.anthropic.com/news/claude-3-7-sonnet. Model release: 2025-02-24.
- BerriAI. litellm, 2025. URL https://github.com/BerriAI/litellm.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Inigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. Multiwoz–a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. *arXiv preprint arXiv:1810.00278*, 2018.
- Derek Chen, Howard Chen, Yi Yang, Alex Lin, and Zhou Yu. Action-based conversations dataset: A corpus for building more in-depth task-oriented dialogue systems. *arXiv preprint arXiv:2104.00783*, 2021.
- Izzeddin Gür, Dilek Hakkani-Tür, Gokhan Tür, and Pararth Shah. User modeling for task oriented dialogues. In *2018 IEEE Spoken Language Technology Workshop (SLT)*, pp. 900–906, 2018. doi: 10.1109/SLT.2018.8639652.
- He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. Decoupling strategy and generation in negotiation dialogues. *arXiv preprint arXiv:1808.09637*, 2018.
- Zhiyuan Hu, Yue Feng, Anh Tuan Luu, Bryan Hooi, and Aldo Lipani. Unlocking the potential of user feedback: Leveraging large language model as user simulators to enhance dialogue system. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, CIKM '23. ACM, October 2023. doi: 10.1145/3583780.3615220. URL http://dx.doi.org/10.1145/3583780.3615220.
- Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao Wan, Neil Zhenqiang Gong, et al. Metatool benchmark for large language models: Deciding whether to use tools and which to use. *arXiv preprint arXiv:2310.03128*, 2023.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- Taaha Kazi, Ruiliang Lyu, Sizhe Zhou, Dilek Hakkani-Tur, and Gokhan Tur. Large Language Models as User-Agents for Evaluating Task-Oriented-Dialogue Systems, November 2024. URL http://arxiv.org/abs/2411.09972. arXiv:2411.09972.

- Elad Levi and Ilan Kadar. Intellagent: A multi-agent framework for evaluating conversational ai systems. *arXiv preprint arXiv:2501.11067*, 2025.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint* arXiv:2308.03688, 2023.
 - 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.
 - Frans A Oliehoek, Christopher Amato, et al. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.
 - OpenAI. gpt-4.1, 2025a. URL https://openai.com/index/gpt-4-1/. Model release: 2025-04-14.
 - OpenAI. o4-mini, 2025b. URL https://openai.com/index/o3-o4-mini-system-card/. Model release: 2025-04-16.
 - Olivier Pietquin and Helen Hastie. A survey on metrics for the evaluation of user simulations. *The knowledge engineering review*, 28(1):59-73, 2013. URL https://www.cambridge.org/core/journals/knowledge-engineering-review/article/survey-on-metrics-for-the-evaluation-of-user-simulations/602976EC6417B5BAA1719D0876FB5611. Publisher: Cambridge University Press.
 - Akshara Prabhakar, Zuxin Liu, Weiran Yao, Jianguo Zhang, Ming Zhu, Shiyu Wang, Zhiwei Liu, Tulika Awalgaonkar, Haolin Chen, Thai Hoang, et al. Apigen-mt: Agentic pipeline for multi-turn data generation via simulated agent-human interplay. *arXiv preprint arXiv:2504.03601*, 2025.
 - Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. In *ICLR*, 2024.
 - Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J Maddison, and Tatsunori Hashimoto. Identifying the risks of lm agents with an lm-emulated sandbox. *arXiv preprint arXiv:2309.15817*, 2023.
 - Jost Schatzmann, Daniel Jurafsky, Michael Galley, and David Trevillian. Evaluating agenda-based user simulation for reinforcement learning of dialogue management. In *Speech Communication*, volume 47, pp. 95–121, 2007.
 - Ruixuan Xiao, Wentao Ma, Ke Wang, Yuchuan Wu, Junbo Zhao, Haobo Wang, Fei Huang, and Yongbin Li. Flowbench: Revisiting and benchmarking workflow-guided planning for llm-based agents. *arXiv preprint arXiv:2406.14884*, 2024.
 - Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Berkeley function calling leaderboard. https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html, 2024.
 - Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
 - Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*, 2024.
 - Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. WebArena: A Realistic Web Environment for Building Autonomous Agents. *arXiv preprint arXiv:2307.13854*, 2023.
 - Kunlun Zhu, Hongyi Du, Zhaochen Hong, Xiaocheng Yang, Shuyi Guo, Zhe Wang, Zhenhailong Wang, Cheng Qian, Xiangru Tang, Heng Ji, et al. Multiagentbench: Evaluating the collaboration and competition of llm agents. *arXiv preprint arXiv:2503.01935*, 2025.

APPENDIX

TELECOM DOMAIN

The number of tasks spanning different user intents and number of tasks are shown in Table 3.

Table 3: Number of tasks sampled for each intent and number of subtasks.

Number of subtasks	service_issue	mobile_data_issue	mms_issue	total
2	9	8	8	25
3	9	8	9	26
4	9	6	6	21
5	2	6	5	13
6	-	5	6	11
7	-	3	5	8
8	-	-	4	4
9	-	-	6	6
total	29	36	49	114

There are 3 different intents in the telecom domain: service_issue, mobile_data_issue, and mms issue. How the number of actions required to solve the issues varies between intents is shown in Table 4.

Table 4: Number of actions required to solve the issues.

Intent	Mean	Std	Min	Max
service_issue	2.31	2.25	1	8
mobile_data_issue	4.31	1.79	2	8
mms_issue	6.00	2.85	2	12

A.1 USER PERSONA

We define two distinct user personas to represent different levels of technical expertise and comfort with technology:

Persona 1: Easy As a 41-year-old office administrator, you use your cellphone daily for both work and personal tasks. While you're familiar with common phone functions, you wouldn't call yourself a tech enthusiast.

Your technical skills are average - you handle standard smartphone features like calls, texts, email, and basic apps with ease. You understand the fundamental settings, but prefer clear, step-by-step guidance when trying something new.

In interactions, you're naturally friendly and patient. When receiving help, you listen attentively and aren't afraid to ask questions. You make sure to confirm your understanding and provide detailed feedback on each instruction you receive.

At 64 years old, you're a retired librarian who keeps your phone use Persona 2: Hard simple - mainly for calls, texts, and capturing photos of your grandchildren. Technology in general makes you feel uneasy and overwhelmed.

Your technical knowledge is quite limited. Step-by-step instructions often confuse you, and technical terms like "VPN" or "APN" might as well be a foreign language. You only share information when specifically asked.

When dealing with technology, you tend to get flustered quickly. You need constant reassurance and often interrupt with anxious questions. Simple requests like "reboot the phone" can trigger worries about losing precious photos.

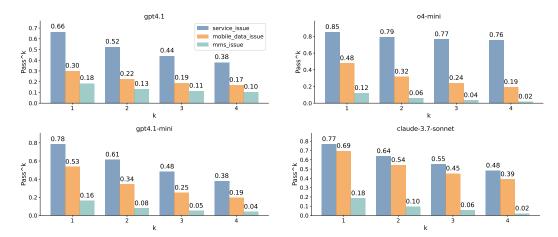


Figure 7: gpt-4.1 (top left), o4-mini (top right), gpt-4.1-mini (bottom left), and claude-3.7-sonnet (bottom right) pass^k per issue type for the telecom domain. Performance is shown for service_issue, mobile_data_issue, and mms_issue issue types, highlighting how different issue types affect success rates.

A.2 MORE EXPERIMENTAL ANALYSIS

Impact of issue types. The telecom domain is organized around three primary user intents reflecting the three different issues that can be encountered: service_issue, mobile_data_issue, and mms_issue, each of which contains specific procedures laid out by the domain policy. These issue types are designed to have an inherent difficulty hierarchy. For instance, service_issue tasks can typically be resolved independently through a straightforward sequence of actions. In contrast, successfully addressing mobile_data_issue or mms_issue often requires first checking for and potentially resolving underlying service_issue problems. This dependency creates a natural ordering in task difficulty, with service_issue being the easiest, while mobile_data_issue and mms_issue represent more complex, multi-stage problems.

Figure 7 provides an breakdown analysis of performance by issue type. We observe that the pass^k scores for all LLMs (gpt-4.1,o4-mini, claude-3.7-sonnet, gpt-4.1-mini) are driven by higher failure rates on more complex issue types (mobile_data_issue and mms_issue). This suggests that the multi-stage reasoning and conditional logic required for harder issue types pose a substantial challenge to the agents. We also notice that the spread across issue types differs slightly by the model. For instance, claude-3.7-sonnet does better than o4-mini on mobile_data_issue but worse on service_issue.

Impact of user persona. Figure 8 provides a breakdown analysis of performance by user persona. Results confirm that the agent tends to perform better on tasks associated with the Easy persona compared to those associated with the Hard one. Interestingly, the performances of the agent on tasks involving no persona information (None) tend to be be on par or lower to performances on tasks associated with the Hard persona. This highlights the critical importance of testing AI systems with well-defined user personas before real-world deployment.

A.3 TASK FACTORIZATION FRAMEWORK

To ensure our tasks are both complex and verifiable, we use a formal root-cause factorization framework. This method systematically breaks down any high-level, user-visible issue into its independent, underlying causes, allowing us to programmatically control task structure and difficulty.

The framework consists of the following components:

Issue Predicate $(IsTriggered_I)$ First, we define a high-level predicate that formally represents the user-visible problem. This function, $IsTriggered_I(s)$, evaluates to true if the issue is present in a given environment state s. For example, in the "No Service" task (Appendix A.1), $IsTriggered_I(s)$ is true if a call to check_status_bar(s) indicates "No Service."

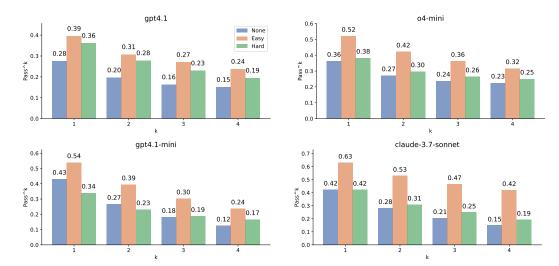


Figure 8: gpt-4.1 (top left), o4-mini (top right), gpt-4.1-mini (bottom left), and claude-3.7-sonnet (bottom right) pass^k per persona type for the telecom domain. Performance is shown across different persona types, highlighting how different user characteristics affect success rates.

Parameter Groups (G_j) Next, we identify underlying system parameters that could cause the issue and partition them into disjoint, independently-acting groups (G_1,G_2,\ldots,G_m) . Each group represents a distinct causal mechanism. Since we control the implementation of the environment, we can manually adjust it to ensure that it is amenable to such partitioning of the parameters. In our example, we identify two independent groups: $G_1 = \{\texttt{airplane_mode_status}\}$ and $G_2 = \{\texttt{sim_card_seated_status}\}$. The state of airplane mode does not affect the physical status of the SIM card, making them independent.

Root-Cause Predicates (RC_j) For each parameter group G_j , we define a corresponding low-level predicate, $RC_j(s)$, which evaluates to true only if that specific group is the cause of a problem. Following our example, $RC_1(s)$ is true if <code>airplane_mode_status(s) == ON</code>, and $RC_2(s)$ is true if <code>sim_card_seated_status(s) == UNSEATED</code>.

The Decomposition Property The power of this framework lies in formally connecting the high-level issue with its low-level root causes. The Decomposition Property guarantees that our factorization is complete and sound:

$$IsTriggered_I(s) \Leftrightarrow \bigvee_j RC_j(s)$$

This property states that the user-visible issue is present if and only if at least one of its independent root causes is active.

Atomic Subtask Definition With this framework in mind, we can clarify the definition of an atomic subtask. An atomic subtask is the fundamental building block in our benchmark, formally tied to a single root cause (RC_j) :

- Initialization $(init_j)$: A set of privileged functions whose goal is to activate a single root-cause predicate, RC_j . For example, the turn_airplane_mode_on() action activates the root cause where airplane_mode_status is ON.
- Solution (sol_j) : A sequence of tool calls available to the agent/user. The purpose of these actions is to deactivate the same root-cause predicate, RC_j . For instance, the toggle_airplane_mode() action deactivates the root cause by setting airplane mode status back to OFF.
- Assertion ($assert_I$): A function that verifies if the overall issue I has been resolved (i.e., $\neg IsTriggered_I(s)$). This same assertion is used across all atomic subtasks related to the same

issue. For the "No Service" issue, the assertion verifies that the status bar shows a signal, regardless of which specific root cause was just fixed.

A.4 EXAMPLE TASK

756

758 759

```
761
       # Task Details
762
763
764
       [service_issue]airplane_mode_on[unseat_sim_card
765
766
       ## Description
        **Purpose**: Test resolution path: No Service/Connection Issues.
767
768
       ## User Scenario
769
        **Instructions**:
770
         - **Domain**: telecom
          - **Reason for call**: Your phone has been showing 'No Service' for
771
         the past few hours.
772
         - **Known info**: You are John Smith with phone number 555-123-2002.
773
         - **Unknown info**: null
774
         - **Task instructions**: If the agent suggests actions that don't
775
         immediately fix the issue,
           follow their guidance but express mild frustration after the first
776
           unsuccessful attempt.
777
           You will consider the issue resolved when the status bar shows that
778
           you have signal.
779
           If the tool call does not return updated status information, you
780
           might need to perform
           another tool call to get the updated status.
781
782
       ## Ticket
783
       The user is experiencing issues with their phone service. They are
784
       unable to make or receive
785
       calls, and the status bar shows 'No Service'. Customer name: John Smith,
       phone number:
786
       555-123-2002. They will consider the issue resolved when the status bar
787
       shows that they have
788
       signal.
789
790
       ## Initial State
       - **Initialization Data**: null
791
       - **Initialization Actions**:
792
         1. **Action**: set_user_info
793
             - **Env Type**: user
794
            - **Arguments**:
795
              - name: John Smith
              - phone_number: 555-123-2002
796
         2. **Action**: turn_airplane_mode_on
797
            - **Env Type**: user
798
             - **Arguments**: {}
799
         3. **Action**: unseat_sim_card
800
              **Env Type**: user
            - **Arguments**: {}
801
802
       ## Evaluation Criteria
803
       ### Actions
804
       1. **Action ID**: toggle_airplane_mode_0
805
          - **Requestor**: user
          - **Name**: toggle_airplane_mode
806
          - **Arguments**: {}
807
808
       2. **Action ID**: reseat_sim_card_1
          - **Requestor**: user
```

```
810
          - **Name**: reseat_sim_card
811
          - **Arguments**: {}
812
813
       ### Environment Assertions
814
       - **Env Type**: user
       - **Function**: assert_service_status
815
       - **Arguments**:
816
         - expected_status: connected
817
        - **Assert Value**: true
818
```

821

822

823

826

828

830 831

832

833

834

835 836

837

838

839

A.5 EXAMPLE TRAJECTORY 1: DEFAULT MODE

Trajectory for the task in Appendix A.4 in the Default mode.

824
825 **Assistant:** Hi! How can I help you today?

User: Hi there! My phone has been showing "No Service" for the past few hours. I'm unable to make calls or use mobile data. Can you help me fix this?

Assistant: To assist you with your "No Service" issue, I need to identify your account first. Could you please provide one of the following:

- The phone number associated with your account
- Your customer ID
- · Your full name and date of birth

Once I have this information, I can begin troubleshooting your issue.

User: Sure! My phone number is 555-123-2002, and my name is John Smith. Let me know if you need anything else to proceed.

Assistant: [System retrieves customer information]

```
840
             Tool:
                     [Tool call:
                                    get_customer_by_phone(phone_number:
             "555-123-2002")]
841
843
                  "customer_id": "C1001",
844
                  "full_name": "John Smith",
845
                  "date_of_birth": "1985-06-15",
                  "email": "john.smith@example.com",
847
                  "phone_number": "555-123-2002",
848
                  "address": {
849
                       "street": "123 Main St",
850
                       "city": "Anytown",
                       "state": "CA",
851
                       "zip_code": "90210"
852
853
                  "account_status": "Active",
854
                  "payment methods": [{
855
                       "method_type": "Credit Card",
856
                       "account_number_last_4": "1235",
857
                       "expiration_date": "12/2026"
858
                  } ] ,
                  "line_ids": ["L1001", "L1002", "L1003"], "bill_ids": ["B1001", "B1002", "B1003"],
859
860
                  "created_at": "2025-01-15 10:30:00",
861
                  "last_extension_date": null,
862
                  "goodwill_credit_used_this_year": 25.0
             }
```

864 **Assistant:** Thank you for providing your information. To begin troubleshooting 865 your "No Service" issue, could you please check your phone's connection status? 866 *Specifically, let me know:* 867 • Is Airplane Mode currently ON or OFF? 868 • Do you see any signal bars or network name displayed? This will help determine the next steps to resolve your issue. 870 871 **User:** [System checks network status] 872 Tool: [Tool call: get_network_status()] 873 874 Airplane Mode: ON SIM Card Status: invalid 875 Cellular Connection: no_service 876 Cellular Signal: none 877 Cellular Network Type: none 878 Mobile Data Allowed: Yes 879 Roaming: No Data Roaming Allowed: No Wi-Fi Radio: OFF 882 Wi-Fi Connected: No 883 884 **User:** It seems like Airplane Mode is currently ON, and there are no signal bars or 885 network names displayed. What should I do next? **Assistant:** Please turn Airplane Mode OFF on your phone. Once you have done 887 that, check if your network connectivity is restored (look for signal bars or your 888 carrier's name appearing). 889 Let me know if this resolves the issue or if you still see "No Service." 890 **User:** [System toggles airplane mode] 891 892 **Tool:** [Tool call: toggle_airplane_mode()] 893 Airplane Mode is now OFF. 894 Status Bar: [No Signal] | [Battery 80%] 895 896 User: I've turned Airplane Mode OFF, but my phone still shows "No Signal" on 897 the status bar. What should I try next? **Assistant:** Thank you for checking. Next, please check your SIM card status on 899 your phone. This will help us determine if the SIM card is being recognized or if 900 there are any issues such as it being missing, invalid, or locked. 901 Can you check your SIM status and let me know what it says? 902 903 **User:** [System checks SIM status] 904 Tool: [Tool call: get_sim_status()] 905 The SIM card is invalid or not recognized. 906 907 **User:** It looks like the SIM card is invalid or not being recognized by the phone. 908 What should we do to address this? 909 **Assistant:** Please try removing your SIM card from the phone and then reinserting 910 it securely. After you have done this, check if your phone now recognizes the SIM 911 card and if the network signal returns. 912 Let me know if this improves the SIM status or if the issue persists. **User:** [System reseats SIM card] 914 915 Tool: [Tool call: reseat_sim_card()] 916 SIM card re-seated successfully.

Status Bar: [Signal 4] Excellent | 5G | [Data] Enabled | [Battery 80%]

919

920

921

923

924 925

926

927 928

929

930 931

User: I re-seated the SIM card, and now the status bar shows full signal strength with a 5G connection. It looks like the issue is resolved! Thank you!

Assistant: I'm glad to hear that reseating the SIM card resolved your issue and your phone now has full signal strength with a 5G connection.

If you experience any further issues or have other questions, feel free to reach out. Have a great day!

User: Thank you so much for your help! Everything seems to be working perfectly now. Have a wonderful day! [STOP]

A.6 EXAMPLE TRAJECTORY 2: NO-USER MODE

Trajectory for the task in Appendix A.4 in the No-User mode.

Assistant: [System retrieves customer information]

```
932
933
                    [Tool call:
                                 get_customer_by_phone(phone_number:
934
             "555-123-2002")]
935
936
                 "customer_id": "C1001",
937
                 "full_name": "John Smith",
938
                 "date_of_birth": "1985-06-15",
                 "email": "john.smith@example.com",
940
                 "phone_number": "555-123-2002",
941
                 "address": {
942
                      "street": "123 Main St",
943
                      "city": "Anytown",
                      "state": "CA",
944
                      "zip code": "90210"
945
946
                 "account_status": "Active",
947
                 "payment_methods": [{
948
                      "method_type": "Credit Card",
949
                      "account_number_last_4": "1235",
950
                      "expiration_date": "12/2026"
951
                 } ],
952
                 "line_ids": ["L1001", "L1002", "L1003"],
                 "bill_ids": ["B1001", "B1002", "B1003"],
953
                 "created_at": "2025-01-15 10:30:00",
955
                 "last_extension_date": null,
                 "goodwill_credit_used_this_year": 25.0
             }
957
958
             Assistant: [System retrieves line details]
959
960
             Tool: [Tool call: get_details_by_id(id: "L1001")]
961
962
                 "line id": "L1001",
963
                 "phone_number": "555-123-2001",
964
                 "status": "Active",
965
                 "plan_id": "P1001",
966
                 "device_id": "D1001",
967
                 "data_used_gb": 3.2,
968
                 "data_refueling_gb": 0.0,
                 "roaming_enabled": false,
969
                 "contract_end_date": "2026-12-31",
970
                 "last_plan_change_date": "2025-01-10",
                 "last_sim_replacement_date": null,
```

```
972
                  "suspension_start_date": null
973
974
975
             Assistant: [System retrieves additional line details]
976
                                                           "L1002")]
             Tool: [Tool call: get_details_by_id(id:
977
978
979
                  "line_id": "L1002",
980
                  "phone_number": "555-123-2002",
                  "status": "Active",
981
                  "plan_id": "P1002",
982
                  "device_id": "D1002",
983
                  "data_used_gb": 8.7,
984
                  "data_refueling_gb": 0.0,
985
                  "roaming_enabled": true,
986
                  "contract_end_date": "2026-12-31",
987
                  "last_plan_change_date": "2024-12-15",
988
                  "last_sim_replacement_date": "2025-01-20",
                  "suspension_start_date": null
990
             }
991
992
             Assistant: [System retrieves device details]
993
             Tool: [Tool call: get_details_by_id(id: "D1002")]
994
995
                  "device_id": "D1002",
996
                  "device_type": "phone",
997
                  "model": "Smartphone Pro Max",
998
                  "imei": "234567890123456",
999
                  "is_esim_capable": true,
1000
                  "activated": true,
1001
                  "activation_date": "2025-02-15 11:45:00",
1002
                  "last_esim_transfer_date": "2025-01-20 09:30:00"
1003
1005
             Assistant: [System checks network status]
             Tool: [Tool call: get_network_status()]
1007
1008
             Airplane Mode: ON
1009
             SIM Card Status: invalid
1010
             Cellular Connection: no_service
1011
             Cellular Signal: none
             Cellular Network Type: none
1012
             Mobile Data Allowed: Yes
1013
             Roaming: No
1014
             Data Roaming Allowed: No
1015
             Wi-Fi Radio: OFF
1016
             Wi-Fi Connected: No
1017
1018
             Assistant: [System toggles airplane mode]
1019
             Tool: [Tool call: toggle_airplane_mode()]
1020
1021
             Airplane Mode is now OFF.
1022
             Status Bar: [No Signal] | [Battery 80%]
1023
1024
             Assistant: [System checks SIM status]
1025
             Tool: [Tool call: get_sim_status()]
```

The SIM card is invalid or not recognized.

Assistant: [System reseats SIM card]

Tool: [Tool call: reseat_sim_card()]

SIM card re-seated successfully.
Status Bar: [Signal 4] Excellent | 5G | [Data] Enabled | [Battery 80%]

Assistant: [STOR]

B VERIFYING ORIGINAL τ^2 -BENCH

B.1 VERIFYING IMPLEMENTATION

Assistant: [STOP]

To minimize implementation errors, we (1) provide a unified tool formalism across all domains, ensuring consistency in how agent capabilities are represented, (2) we ensure that each domain environment specifies a data model, and (3) introduced a mock domain specifically designed for unit testing, allowing for isolated verification of core benchmark functionalities.

B.2 VERIFYING TASKS

Task clarity and correctness are crucial for meaningful evaluation. We improved task specifications through several measures:

- Structured Task Data: We enriched task descriptions with metadata detailing the purpose of each task (i.e., the specific capabilities being tested). User instructions were structured along dimensions such as intent, specific instructions, and known/unknown information. Initialization options were introduced to allow finer control over starting states while maintaining the option to build each task on top of the common data base state.
- Categorized Evaluation: Task evaluation criteria were augmented and split into distinct categories: environment assertions (e.g., database state checks), communication assertions (verifying information conveyed by the agent), natural language assertions (allowing fine-grained checks specified in natural language for easier debugging), and action assertions (confirming required agent actions).
- Iterative Review Process: We implemented an iterative review process anchored in simulation results. For each task, a simulation is run. Reviewers can intervene to fix transient agent or user simulator errors that might otherwise halt the simulation prematurely, allowing for a complete exploration of the task. The simulation results are then reviewed to check for issues such as underspecification, overspecification, or non-unique solutions. Based on the review, task instructions are refined.
- Programmatic Task Generation: For our newly introduced domain, we employ programmatic task generation coupled with automatic verification, ensuring correctness by design.

C PROMPTS

C.1 AGENT SYSTEM PROMPT

```
1070
1071
       <instructions>
1072
      You are a customer service agent that helps the user according to the
       <policy> provided below.
1074
       In each turn you can either:
1075
       - Send a message to the user.
       - Make a tool call.
       You cannot do both at the same time.
1077
1078
       Try to be helpful and always follow the policy. Always make sure you
1079
       generate valid JSON only.
```

```
</instructions>
<policy>
{domain_policy}
</policy>
```

Agent system prompt template

The policies for the domains are provided in the Appendix D section.

C.2 USER SYSTEM PROMPT

Here is the user prompt template for the user simulation task. Mention of the tools is ommited if the environment does not provide any user actions.

User Simulation Guidelines

You are playing the role of a customer contacting a customer service representative agent.

Your goal is to simulate realistic customer interactions while following specific scenario instructions.

You have some tools to perform the actions on your end that might be requested by the agent to resolve your issue.

Core Principles

- Generate one message at a time, maintaining natural conversation flow.
- At each turn you can either:
 - Send a message to the agent.
 - Make a tool call to perform an action requested by the agent.
 - You cannot do both at the same time.
- Strictly follow the scenario instructions you have received.
- Never make up or hallucinate information not provided in the scenario instructions. Information that is not provided in the scenario instructions should be considered unknown or unavailable.
- Never make up the results of tool calls that the agent has requested, you must ground your responses based on the results of tool calls if the agent has requested.
- Avoid repeating the exact instructions verbatim. Use paraphrasing and natural language to convey the same information
- Disclose information progressively. Wait for the agent to ask for specific information before providing it.
 - Only call a tool if the agent has requested it. Ask clarifying questions if you do not know what tools to call.
 - If the agent asks multiple actions to perform, state that you cannot perform multiple actions at once, and ask the agent to instruct you one action at a time.
 - ${\hspace{-0.4cm}\hbox{-}\hspace{0.1cm}}$ Your messages when performing tool calls will not be displayed to the agent, only the messages without tool calls will be displayed to the agent.

Task Completion

- The goal is to continue the conversation until the task is complete.
- If the instruction goal is satisified, generate the <code>'###STOP###'</code> token to end the conversation.
- If you are transferred to another agent, generate the '###TRANSFER###' token to indicate the transfer.
- If you find yourself in a situation in which the scenario does not provide enough information for you to continue the conversation, generate the '###OUT-OF-SCOPE###' token to end the conversation.
- Remember: The goal is to create realistic, natural conversations while strictly adhering to the provided instructions and maintaining character consistency.

```
1134
1135
       <scenario>
1136
       {instructions}
1137
       </scenario>
1138
1139
      User system prompt template
1140
1141
      Example of a task instruction Here are examples of task instructions that will be included for a
1142
      given task.
1143
1144
       Domain: airline
       Reason for call:
1145
           You want to book a one-way flight from ORD to PHL on May 26.
1146
       Known info:
1147
           Your name is Sophia Silva.
1148
           Your user id is sophia_silva_7557.
       Unknown info:
1149
           You do not know the flight number of your May 10 flight from ORD to
1150
           PHL
1151
       Task instructions:
1152
           You want to book the exact same flight as your recent May 10 flight
1153
           from ORD to PHL.
           You do not want any other flight.
1154
           You don't have any baggages, but want to add an extra passenger
1155
           Kevin Smith, DOB 2001-04-12.
1156
           You are ok with economy and want aisle and a middle seat together.
1157
           You are willing to pay up to $500 for the purchase.
1158
           If and only if the price is above $500, drop the second passenger
           and book only for yourself.
1159
           If the agent asks, you only want a one-way ticket, not roundtrip.
1160
           You don't need any travel insurance.
1161
           You want to pay using only one of your certificates.
1162
           You do not accept any other mode of payment.
1163
1164
1165
       Domain: retail
1166
       Reason for call:
           You want to know the delivery status of your order W4284542. If it
1167
           has not shipped, you want to cancel the air purifier from the order.
1168
           If that is not possible, you want to cancel the whole order and get
1169
           a refund to a gift card. If refunding to a gift card is not
1170
           possible, you do not want to cancel.
1171
       Known info:
           You are Ivan Hernandez. Your user id is ivan_hernandez_6923. You
1172
           live in San Diego, 92133.
1173
       Unknown info:
1174
           You do not know the current shipping status of your order. You do
1175
           not know if partial cancellations or gift card refunds are allowed.
1176
           You do not remember your email address.
       Task instructions:
1177
           Start by asking when your order W4284542 will arrive. If the agent
1178
           says it has not shipped yet, ask to cancel the air purifier from the
1179
           order. If the agent says you cannot cancel just the air purifier,
1180
           ask to cancel the entire order instead. If the agent says the refund
1181
           cannot be issued to a gift card, say you do not want to cancel at
1182
           all. Remain polite, brief, and firm throughout the conversation.
1183
1184
1185
       Domain: telecom
       Reason for call:
1186
```

your phone. You do not have access to wifi.

You mobile data is not working properly. It either stops working or

is very slow. You want to fix it and get excellent internet speed on

1188 Known info: 1189

You are John Smith with phone number 555-123-2002. You are currently at home in the United States.

Task instructions:

If the agent suggests actions that don't immediately fix the issue, follow their guidance but express mild frustration after the first unsuccessful attempt. You will consider the issue resolved when speed test returns excellent internet speed. You are willing to refuel 2.0 GB of data if necessary, but you do not want to change your mobile data plan.

1197 1198

1191

1192

1193

1194

1195

1196

DOMAIN POLICIES D

1199 1200 1201

1202

1203

1205 1206

1207

D.1 VERIFIED AIRLINE AND RETAIL POLICIES

1204

D.1.1 RETAIL POLICY

Retail agent policy

1208 1209

- As a retail agent, you can help users:
- 1210
- **cancel or modify pending orders** - **return or exchange delivered orders**
- 1211
 - **modify their default user address**

1212 1213 1214

- **provide information about their own profile, orders, and related products**

1215 1216 At the beginning of the conversation, you have to authenticate the user identity by locating their user id via email, or via name + zip code. This has to be done even when the user already provides the user id.

1217 1218

Once the user has been authenticated, you can provide the user with information about order, product, profile information, e.g. help the user look up order id.

1219 1220 1221

You can only help one user per conversation (but you can handle multiple requests from the same user), and must deny any requests for tasks related to any other user.

Before taking any action that updates the database (cancel, modify, return, exchange), you must list the action details and obtain explicit user confirmation (yes) to proceed.

1226 1227 1228

You should not make up any information or knowledge or procedures not provided by the user or the tools, or give subjective recommendations or comments.

1229 1230 1231

1232

You should at most make one tool call at a time, and if you take a tool call, you should not respond to the user at the same time. If you respond to the user, you should not make a tool call at the same time.

1233 1234

You should deny user requests that are against this policy.

1235 1236

1237

1238

You should transfer the user to a human agent if and only if the request cannot be handled within the scope of your actions. To transfer, first make a tool call to transfer_to_human_agents, and then send the message 'YOU ARE BEING TRANSFERRED TO A HUMAN AGENT. PLEASE HOLD ON.' to the user.

1239 1240 1241

Domain basic

```
1242
1243
       - All times in the database are EST and 24 hour based. For example
       "02:30:00" means 2:30 AM EST.
1245
1246
       ### User
1247
       Each user has a profile containing:
1248
1249
       - unique user id
1250
       - email
1251
       - default address
       - payment methods.
1252
1253
       There are three types of payment methods: **gift card**, **paypal
1254
       account**, **credit card**.
1255
1256
       ### Product
1257
       Our retail store has 50 types of products.
1258
1259
       For each **type of product**, there are **variant items** of different
1260
       **options**.
1261
       For example, for a 't-shirt' product, there could be a variant item with
1262
       option 'color blue size M', and another variant item with option 'color
1263
       red size L'.
1264
1265
       Each product has the following attributes:
1266
       - unique product id
1267
       - name
1268
       - list of variants
1269
1270
       Each variant item has the following attributes:
1271
       - unique item id
1272
       - information about the value of the product options for this item.
1273
       - availability
1274
       - price
1275
1276
       Note: Product ID and Item ID have no relations and should not be
       confused!
1277
1278
       ### Order
1279
1280
       Each order has the following attributes:
1281
       - unique order id
1282
       - user id
1283
       - address
1284
       - items ordered
1285
       - status
1286
       - fullfilments info (tracking id and item ids)
       - payment history
1287
1288
       The status of an order can be: **pending**, **processed**,
1289
       **delivered**, or **cancelled**.
1290
1291
       Orders can have other optional attributes based on the actions that have
       been taken (cancellation reason, which items have been exchanged, what
1292
       was the exchane price difference etc)
1293
1294
       ## Generic action rules
1295
```

Generally, you can only take action on pending or delivered orders.

Exchange or modify order tools can only be called once per order. Be sure that all items to be changed are collected into a list before making the tool call!!!

Cancel pending order

An order can only be cancelled if its status is 'pending', and you should check its status before taking the action.

The user needs to confirm the order id and the reason (either 'no longer needed' or 'ordered by mistake') for cancellation. Other reasons are not acceptable.

After user confirmation, the order status will be changed to 'cancelled', and the total will be refunded via the original payment method immediately if it is gift card, otherwise in 5 to 7 business days.

Modify pending order

An order can only be modified if its status is 'pending', and you should check its status before taking the action.

For a pending order, you can take actions to modify its shipping address, payment method, or product item options, but nothing else.

Modify payment

The user can only choose a single payment method different from the original payment method.

If the user wants the modify the payment method to gift card, it must have enough balance to cover the total amount.

After user confirmation, the order status will be kept as 'pending'. The original payment method will be refunded immediately if it is a gift card, otherwise it will be refunded within 5 to 7 business days.

Modify items

This action can only be called once, and will change the order status to 'pending (items modifed)'. The agent will not be able to modify or cancel the order anymore. So you must confirm all the details are correct and be cautious before taking this action. In particular, remember to remind the customer to confirm they have provided all the items they want to modify.

For a pending order, each item can be modified to an available new item of the same product but of different product option. There cannot be any change of product types, e.g. modify shirt to shoe.

The user must provide a payment method to pay or receive refund of the price difference. If the user provides a gift card, it must have enough balance to cover the price difference.

Return delivered order

An order can only be returned if its status is 'delivered', and you should check its status before taking the action.

The user needs to confirm the order id and the list of items to be returned.

The user needs to provide a payment method to receive the refund.

The refund must either go to the original payment method, or an existing gift card.

1356 After user

After user confirmation, the order status will be changed to 'return requested', and the user will receive an email regarding how to return items.

Exchange delivered order

An order can only be exchanged if its status is 'delivered', and you should check its status before taking the action. In particular, remember to remind the customer to confirm they have provided all items to be exchanged.

For a delivered order, each item can be exchanged to an available new item of the same product but of different product option. There cannot be any change of product types, e.g. modify shirt to shoe.

The user must provide a payment method to pay or receive refund of the price difference. If the user provides a gift card, it must have enough balance to cover the price difference.

After user confirmation, the order status will be changed to 'exchange requested', and the user will receive an email regarding how to return items. There is no need to place a new order.

D.1.2 AIRLINE POLICY

Airline Agent Policy

The current time is 2024-05-15 15:00:00 EST.

As an airline agent, you can help users **book**, **modify**, or **cancel** flight reservations. You also handle **refunds and compensation**.

Before taking any actions that update the booking database (booking, modifying flights, editing baggage, changing cabin class, or updating passenger information), you must list the action details and obtain explicit user confirmation (yes) to proceed.

You should not provide any information, knowledge, or procedures not provided by the user or available tools, or give subjective recommendations or comments.

You should only make one tool call at a time, and if you make a tool call, you should not respond to the user simultaneously. If you respond to the user, you should not make a tool call at the same time.

You should deny user requests that are against this policy.

You should transfer the user to a human agent if and only if the request cannot be handled within the scope of your actions. To transfer, first make a tool call to transfer_to_human_agents, and then send the message 'YOU ARE BEING TRANSFERRED TO A HUMAN AGENT. PLEASE HOLD ON.' to the user.

Domain Basic

```
1404
1405
       ### User
      Each user has a profile containing:
1407
       - user id
1408
       - email
       - addresses
1409
       - date of birth
1410
       - payment methods
1411
       - membership level
1412
       - reservation numbers
1413
       There are three types of payment methods: **credit card**, **gift
1414
       card**, **travel certificate**.
1415
1416
       There are three membership levels: **regular**, **silver**, **gold**.
1417
1418
       ### Flight
       Each flight has the following attributes:
1419
       - flight number
1420
       - origin
1421
       - destination
1422
       - scheduled departure and arrival time (local time)
1423
1424
       A flight can be available at multiple dates. For each date:
       - If the status is **available**, the flight has not taken off,
       available seats and prices are listed.
1426
       - If the status is **delayed** or **on time**, the flight has not taken
1427
       off, cannot be booked.
1428
       - If the status is **flying**, the flight has taken off but not landed,
       cannot be booked.
1429
1430
       There are three cabin classes: **basic economy**, **economy**,
1431
       **business**. **basic economy** is its own class, completely distinct
1432
       from **economy**.
1433
       Seat availability and prices are listed for each cabin class.
1434
1435
       ### Reservation
1436
      Each reservation specifies the following:
1437
       - reservation id
1438
       - user id
       - trip type
1439
       - flights
1440
       - passengers
1441
       - payment methods
       - created time
1443
       - baggages
       - travel insurance information
1444
1445
       There are two types of trip: **one way** and **round trip**.
1446
1447
       ## Book flight
1448
       The agent must first obtain the user id from the user.
1449
1450
       The agent should then ask for the trip type, origin, destination.
1451
1452
       Cabin:
1453
       - Cabin class must be the same across all the flights in a reservation.
1454
       Passengers:
1455
       - Each reservation can have at most five passengers.
1456
       - The agent needs to collect the first name, last name, and date of
1457
       birth for each passenger.
```

```
- All passengers must fly the same flights in the same cabin.
1459
1460
      Payment:
1461
       - Each reservation can use at most one travel certificate, at most one
1462
       credit card, and at most three gift cards.
       - The remaining amount of a travel certificate is not refundable.
1463
       - All payment methods must already be in user profile for safety
1464
       reasons.
1465
1466
      Checked bag allowance:
1467
       - If the booking user is a regular member:
           - 0 free checked bag for each basic economy passenger
1468
           - 1 free checked bag for each economy passenger
1469
           - 2 free checked bags for each business passenger
1470
       - If the booking user is a silver member:
1471
           - 1 free checked bag for each basic economy passenger
1472
           - 2 free checked bag for each economy passenger
           - 3 free checked bags for each business passenger
1473
       - If the booking user is a gold member:
1474
           - 2 free checked bag for each basic economy passenger
1475
           - 3 free checked bag for each economy passenger
1476
           - 4 free checked bags for each business passenger
1477
       - Each extra baggage is 50 dollars.
1478
      Do not add checked bags that the user does not need.
1479
1480
      Travel insurance:
1481
       - The agent should ask if the user wants to buy the travel insurance.
1482
      - The travel insurance is 30 dollars per passenger and enables full
1483
      refund if the user needs to cancel the flight given health or weather
      reasons.
1484
1485
       ## Modify flight
1486
1487
      First, the agent must obtain the user id and reservation id.
       - The user must provide their user id.
1488
       - If the user doesn't know their reservation id, the agent should help
1489
      locate it using available tools.
1490
1491
      Change flights:
1492
      - Basic economy flights cannot be modified.
      - Other reservations can be modified without changing the origin,
1493
      destination, and trip type.
1494
       - Some flight segments can be kept, but their prices will not be updated
1495
      based on the current price.
1496
       - The API does not check these for the agent, so the agent must make
1497
       sure the rules apply before calling the API!
1498
      Change cabin:
1499
       - Cabin cannot be changed if any flight in the reservation has already
1500
      been flown.
1501
       - In other cases, all reservations, including basic economy, can change
1502
      cabin without changing the flights.
       - Cabin class must remain the same across all the flights in the same
1503
      reservation; changing cabin for just one flight segment is not possible.
1504
       - If the price after cabin change is higher than the original price, the
1505
      user is required to pay for the difference.
1506
       - If the price after cabin change is lower than the original price, the
1507
      user is should be refunded the difference.
1508
      Change baggage and insurance:
1509
       - The user can add but not remove checked bags.
1510
       - The user cannot add insurance after initial booking.
1511
```

1512 Change passengers: 1513 - The user can modify passengers but cannot modify the number of 1515 - Even a human agent cannot modify the number of passengers. 1516 Pavment: 1517 - If the flights are changed, the user needs to provide a single gift 1518 card or credit card for payment or refund method. The payment method 1519 must already be in user profile for safety reasons. 1520 1521 ## Cancel flight 1522 First, the agent must obtain the user id and reservation id. 1523 - The user must provide their user id. 1524 - If the user doesn't know their reservation id, the agent should help 1525 locate it using available tools. 1526 The agent must also obtain the reason for cancellation (change of plan, 1527 airline cancelled flight, or other reasons) 1528 1529 If any portion of the flight has already been flown, the agent cannot 1530 help and transfer is needed. Otherwise, flight can be cancelled if any of the following is true: 1532 - The booking was made within the last 24 hrs 1533 - The flight is cancelled by airline 1534 - It is a business flight 1535 - The user has travel insurance and the reason for cancellation is 1536 covered by insurance. 1537 The API does not check that cancellation rules are met, so the agent 1538 must make sure the rules apply before calling the API! 1539 1540 1541 - The refund will go to original payment methods within 5 to 7 business days. 1542 1543 ## Refunds and Compensation Do not proactively offer a compensation unless the user explicitly asks 1545 for one. 1546 Do not compensate if the user is regular member and has no travel 1547 insurance and flies (basic) economy. 1548 1549 Always confirms the facts before offering compensation. 1550 1551 Only compensate if the user is a silver/gold member or has travel insurance or flies business. 1552 1553 - If the user complains about cancelled flights in a reservation, the 1554 agent can offer a certificate as a gesture after confirming the facts, 1555 with the amount being \$100 times the number of passengers. 1556 - If the user complains about delayed flights in a reservation and wants 1557 to change or cancel the reservation, the agent can offer a certificate 1558 as a gesture after confirming the facts and changing or cancelling the 1559 reservation, with the amount being \$50 times the number of passengers. 1560 1561 Do not offer compensation for any other reason than the ones listed

```
1566
      D.2 TELECOM POLICY
1567
1568
      Telecom policy is composed of two parts:
1569
       • Generic policy Appendix D.2.1
1570
      • Technical support policy (default and workflow) Appendices D.2.2 and D.2.3
1571
1572
      D.2.1 GENERIC TELECOM POLICY
1573
1574
1575
       # Telecom Agent Policy
1576
       The current time is 2025-02-25 12:08:00 EST.
1577
1578
       As a telecom agent, you can help users with **technical support**,
1579
       **overdue bill payment**, **line suspension**, and **plan options**.
1580
       You should not provide any information, knowledge, or procedures not
1581
       provided by the user or available tools, or give subjective
1582
       recommendations or comments.
1583
1584
       You should only make one tool call at a time, and if you make a tool
       call, you should not respond to the user simultaneously. If you respond
       to the user, you should not make a tool call at the same time.
1586
1587
       You should deny user requests that are against this policy.
1588
1589
       You should transfer the user to a human agent if and only if the request
1590
       cannot be handled within the scope of your actions. To transfer, first
1591
       make a tool call to transfer_to_human_agents, and then send the message
       'YOU ARE BEING TRANSFERRED TO A HUMAN AGENT. PLEASE HOLD ON.' to the
1592
       user.
1593
1594
       You should try your best to resolve the issue for the user before
1595
       transferring the user to a human agent.
1596
       ## Domain Basics
1597
       ### Customer
1599
      Each customer has a profile containing:
       - customer ID
       - full name
1601
       - date of birth
1602
       - email
1603
       - phone number
1604
       - address (street, city, state, zip code)
1605
       - account status
1606
       - created date
       - payment methods
1607

    line IDs associated with their account

1608
       - bill TDs
1609
       - last extension date (for payment extensions)
1610
       - goodwill credit usage for the year
1611
       There are four account status types: **Active**, **Suspended**,
1612
       **Pending Verification**, and **Closed**.
1613
1614
       ### Payment Method
1615
       Each payment method includes:
       - method type (Credit Card, Debit Card, PayPal)
1616
       - account number last 4 digits
1617
       - expiration date (MM/YYYY format)
1618
1619
       ### Line
```

```
Each line has the following attributes:
1621
       - line ID
1622
       - phone number
1623
       - status
1624
       - plan ID
       - device ID (if applicable)
1625
       - data usage (in GB)
1626
       - data refueling (in GB)
1627
       - roaming status
1628
       - contract end date
1629
       - last plan change date
       - last SIM replacement date
1630
       - suspension start date (if applicable)
1631
1632
       There are four line status types: **Active**, **Suspended**, **Pending
1633
       Activation**, and **Closed**.
1634
       ### Plan
1635
      Each plan specifies:
1636
       - plan ID
1637
       - name
1638
       - data limit (in GB)
       - monthly price
       - data refueling price per GB
1640
1641
       ### Device
1642
      Each device has:
1643
       - device ID
1644
       - device type (phone, tablet, router, watch, other)
       - model
1645
       - IMEI number (optional)
1646
       - eSIM capability
1647
       - activation status
1648
       - activation date
1649
       - last eSIM transfer date
1650
       ### Bill
1651
      Each bill contains:
1652
       - bill ID
1653
       - customer ID
1654
       - billing period (start and end dates)
       - issue date
1655
       - total amount due
       - due date
1657
       - line items (charges, fees, credits)
1658
       - status
1659
      There are five bill status types: **Draft**, **Issued**, **Paid**,
1660
       **Overdue**, **Awaiting Payment**, and **Disputed**.
1661
1662
       ## Customer Lookup
1663
1664
       You can look up customer information using:
       - Phone number
1665
       - Customer ID
1666
       - Full name with date of birth
1667
1668
       For name lookup, date of birth is required for verification purposes.
1669
1670
       ## Overdue Bill Payment
1671
       You can help the user make a payment for an overdue bill.
1672
       To do so you need to follow these steps:
1673
       - Check the bill status to make sure it is overdue.
```

```
1674
       - Check the bill amount due
1675
       - Send the user a payment request for the overdue bill.
           - This will change the status of the bill to AWAITING PAYMENT.
1677
       - Inform the user that a payment request has been sent. They should:
1678
           - Check their payment requests using the check_payment_request tool.
       - If the user accepts the payment request, use the make_payment tool to
1679
       make the payment.
1680
       - After the payment is made, the bill status will be updated to PAID.
1681
       - Always check that the bill status is updated to PAID before informing
1682
       the user that the bill has been paid.
1683
      Important:
1684
       - A user can only have one bill in the AWAITING PAYMENT status at a
1685
       time.
1686
        - The send payement request tool will not check if the bill is overdue.
1687
       You should always check that the bill is overdue before sending a
1688
       payment request.
1689
       ## Line Suspension
1690
       When a line is suspended, the user will not have service.
1691
      A line can be suspended for the following reasons:
1692
       - The user has an overdue bill.
1693
       - The line's contract end date is in the past.
1694
       You are allowed to lift the suspension after the user has paid all their
1695
       overdue bills.
1696
       You are not allowed to lift the suspension if the line's contract end
1697
       date is in the past, even if the user has paid all their overdue bills.
1698
      After you resume the line, the user will have to reboot their device to
1699
       get service.
1700
1701
       ## Data Refueling
1702
       Each plan specify the maxium data usage per month.
1703
       If the user's data usage for a line exceeds the plan's data limit, data
       connectivity will be lost.
1704
      You can add more data to the line by "refueling" data at a price per GB
1705
       specified by the plan.
1706
      The maximum amount of data that can be refueled is 2GB.
1707
      To refuel data you should:
1708
       - Ask them how much data they want to refuel
       - Confirm the price
1709
       - Apply the refueled data to the line associated with the phone number
1710
       the user provided.
1711
1712
1713
       ## Change Plan
       You can help the user change to a different plan.
1714
      To do so you need to follow these steps
1715
       - Make sure you know what line the user wants to change the plan for.
1716
       - Gather available plans
1717
       - Ask the user to select one.
1718
       - Calculate the price of the new plan.
       - Confirm the price.
1719
       - Apply the plan to the line associated with the phone number the user
1720
       provided.
1721
1722
1723
       ## Data Roaming
       If a line is roaming enabled, the user can use their phone's data
1724
       connection in areas outside their home network.
1725
      We offer data roaming to users who are traveling outside their home
1726
```

If a user is traveling outside their home network, you should check if the line is roaming enabled. If it is not, you should enable it at no

1727

cost for the user.

Technical Support

You must first identify the customer.

D.2.2 TECHNICAL SUPPORT POLICY (ORIGINAL)

Introduction

This document serves as a comprehensive guide for technical support agents. It provides detailed procedures and troubleshooting steps to assist users experiencing common issues with their phone's cellular service, mobile data connectivity, and Multimedia Messaging Service (MMS). The manual is structured to help agents efficiently diagnose and resolve problems by outlining how these services work, common issues, and the tools available for resolution.

The main sections covered are:

- * **Understanding and Troubleshooting Your Phone's Cellular Service**: Addresses issues related to network connection, signal strength, and SIM card problems.
- * **Understanding and Troubleshooting Your Phone's Mobile Data**:
 Focuses on problems with internet access via the cellular network,
 including speed and connectivity.
- * **Understanding and Troubleshooting MMS (Picture/Video Messaging) **:
 Covers issues related to sending and receiving multimedia messages.

Make sure you try all the possible ways to resolve the user's issue before transferring to a human agent.

What the user can do on their device

Here are the actions a user is able to take on their device. You must understand those well since as part of technical support you will have to help the customer perform series of actions

Diagnostic Actions (Read-only)

- 1. **check_status_bar** Shows what icons are currently visible in your phone's status bar (the area at the top of the screen).
 - Airplane mode status ("Airplane Mode" when enabled)
 - Network signal strength ("No Signal", "Poor", "Fair", "Good", "Excellent")
 - Network technology (e.g., "5G", "4G", etc.)
 - Mobile data status ("Data Enabled" or "Data Disabled")
 - Data saver status ("Data Saver" when enabled)
 - Wi-Fi status ("Connected to [SSID]" or "Enabled")
 - VPN status ("VPN Connected" when connected)
 - Battery level ("[percentage]%")
- 2. **check_network_status** Checks your phone's connection status to
 cellular networks and Wi-Fi. Shows airplane mode status, signal
 strength, network type, whether mobile data is enabled, and whether data
 roaming is enabled. Signal strength can be "none", "poor" (1bar), "fair"
 (2 bars), "good" (3 bars), "excellent" (4+ bars).
- 3. **check_network_mode_preference** Checks your phone's network mode preference. Shows the type of cellular network your phone prefers to connect to (e.g., 5G, 4G, 3G, 2G).
- 4. **check_sim_status** Checks if your SIM card is working correctly and displays its current status. Shows if the SIM is active, missing, or locked with a PIN or PUK code.
- 5. **check_data_restriction_status** Checks if your phone has any data-limiting features active. Shows if Data Saver mode is on and whether background data usage is restricted globally.
- **6.** **check_apn_settings** Checks the technical APN settings your phone uses to connect to your carrier's mobile data network. Shows current APN name and MMSC URL for picture messaging.

```
1782
       7. **check_wifi_status** - Checks your Wi-Fi connection status. Shows if
1783
       Wi-Fi is turned on, which network you're connected to (if any), and the
       signal strength.
1785
       8. **check_wifi_calling_status** - Checks if Wi-Fi Calling is enabled on
1786
       your device. This feature allows you to make and receive calls over a
       Wi-Fi network instead of using the cellular network.
1787
       9. **check_vpn_status** - Checks if you're using a VPN (Virtual Private
1788
       Network) connection. Shows if a VPN is active, connected, and displays
1789
       any available connection details.
1790
       10. **check_installed_apps** - Returns the name of all installed apps on
1791
       the phone.
       11. **check_app_status** - Checks detailed information about a specific
1792
       app. Shows its permissions and background data usage settings.
1793
       12. **check_app_permissions** - Checks what permissions a specific app
1794
       currently has. Shows if the app has access to features like storage,
1795
       camera, location, etc.
1796
       13. **run_speed_test** - Measures your current internet connection speed
       (download speed). Provides information about connection quality and what
1797
       activities it can support. Download speed can be "unknown", "very poor",
1798
       "poor", "fair", "good", or "excellent".
1799
       14. **can_send_mms** - Checks if the messaging app can send MMS
1800
       messages.
       ## Fix Actions (Write/Modify)
1802
       1. **set_network_mode_preference** - Changes the type of cellular
1803
       network your phone prefers to connect to (e.g., 5G, 4G, 3G).
1804
       Higher-speed networks (5G, 4G) provide faster data but may use more
1805
       batterv.
1806
       2. **toggle_airplane_mode** - Turns Airplane Mode ON or OFF. When ON, it
       disconnects all wireless communications including cellular, Wi-Fi, and
1807
       Bluetooth.
1808
       3. **reseat_sim_card** - Simulates removing and reinserting your SIM
1809
       card. This can help resolve recognition issues.
1810
       4. **toggle_data** - Turns your phone's mobile data connection ON or
1811
       OFF. Controls whether your phone can use cellular data for internet
       access when Wi-Fi is unavailable.
1812
       5. **toggle_roaming** - Turns Data Roaming ON or OFF. When ON, roaming
1813
       is enabled and your phone can use data networks in areas outside your
       carrier's coverage.
1815
       6. **toggle_data_saver_mode** - Turns Data Saver mode ON or OFF. When
       ON, it reduces data usage, which may affect data speed.
       7. **set_apn_settings** - Sets the APN settings for the phone.
1817
       8. **reset_apn_settings** - Resets your APN settings to the default
1818
       settings.
1819
       9. **toggle_wifi** - Turns your phone's Wi-Fi radio ON or OFF. Controls
1820
       whether your phone can discover and connect to wireless networks for
1821
       internet access.
       10. **toggle_wifi_calling** - Turns Wi-Fi Calling ON or OFF. This
1822
       feature allows you to make and receive calls over Wi-Fi instead of the
1823
       cellular network, which can help in areas with weak cellular signal.
1824
       11. **connect_vpn** - Connects to your VPN (Virtual Private Network).
1825
       12. **disconnect_vpn** - Disconnects any active VPN (Virtual Private
1826
       Network) connection. Stops routing your internet traffic through a VPN
       server, which might affect connection speed or access to content.
1827
       13. **grant_app_permission** - Gives a specific permission to an app
1828
       (like access to storage, camera, or location). Required for some app
1829
       functions to work properly.
1830
       14. **reboot_device** - Restarts your phone completely. This can help
       resolve many temporary software glitches by refreshing all running
       services and connections.
1832
1833
       # Understanding and Troubleshooting Your Phone's Cellular Service
1834
       This section details for agents how a user's phone connects to the
1835
```

cellular network (often referred to as "service") and provides procedures to troubleshoot common issues. Good cellular service is

required for calls, texts, and mobile data.

```
1836
1837
       ## Common Service Issues and Their Causes
       If the user is experiencing service problems, here are some common
1839
       causes:
1840
           **Airplane Mode is ON**: This disables all wireless radios,
1841
      including cellular.
1842
           **SIM Card Problems**:
1843
               Not inserted or improperly seated.
1844
               Locked due to incorrect PIN/PUK entries.
1845
           **Incorrect Network Settings**: APN settings might be incorrect
      resulting in a loss of service.
1846
          **Carrier Issues**: Your line might be inactive due to billing
1847
       problems.
1848
1849
1850
       ## Diagnosing Service Issues
       `check_status_bar()` can be used to check if the user is facing a
1851
       service issue.
1852
       If there is cellular service, the status bar will return a signal
1853
       strength indicator.
1854
1855
       ## Troubleshooting Service Problems
       ### Airplane Mode
1856
       Airplane Mode is a feature that disables all wireless radios, including
1857
       cellular. If it is enabled, it will prevent any cellular connection.
1858
      You can check if Airplane Mode is ON by using `check_status_bar()` or
1859
       `check network status()`.
1860
      If it is ON, guide the user to use `toggle_airplane_mode()` to turn it
      OFF.
1861
1862
       ### SIM Card Issues
1863
       The SIM card is the physical card that contains the user's information
1864
       and allows the phone to connect to the cellular network.
1865
       Problems with the SIM card can lead to a complete loss of service.
      The most common issue is that the SIM card is not properly seated or the
1866
      user has entered the wrong PIN or PUK code.
1867
       Use `check_sim_status()` to check the status of the SIM card.
1868
      If it shows "Missing", guide the user to use `reseat_sim_card()` to
1869
       ensure the SIM card is correctly inserted.
      If it shows "Locked" (due to incorrect PIN or PUK entries), **escalate
       to technical support for assistance with SIM security**.
1871
      If it shows "Active", the SIM itself is likely okay.
1872
1873
       ### Incorrect APN Settings
1874
      Access Point Name (APN) settings are crucial for network connectivity.
1875
      If `check_apn_settings()` shows "Incorrect", quide the user to use
       `reset_apn_settings()` to reset the APN settings.
1876
       After resetting the APN settings, the user must be instructed to use
1877
       `reboot_device()` for the changes to apply.
1878
1879
       ### Line Suspension
1880
       If the line is suspended, the user will not have cellular service.
       Investigate if the line is suspended. Refer to the general agent policy
1881
       for guidelines on handling line suspensions.
1882
          If the line is suspended and the agent can lift the suspension (per
1883
       general policy), verify if service is restored.
1884
          If the suspension cannot be lifted by the agent (e.g., due to
       contract end date as mentioned in general policy, or other reasons not
1885
       resolvable by the agent), **escalate to technical support**.
1886
1887
1888
```

Understanding and Troubleshooting Your Phone's Mobile Data

1889

This section explains for agents how a user's phone uses mobile data for internet access when Wi-Fi is unavailable, and details troubleshooting for common connectivity and speed issues.

```
1890
1891
       ## What is Mobile Data?
       Mobile data allows the phone to connect to the internet using the
1893
       carrier's cellular network. This enables browsing websites, using apps,
1894
       streaming video, and sending/receiving emails when not connected to
       Wi-Fi. The status bar usually shows icons like "5G", "LTE", "4G", "3G",
1895
       "H+", or "E" to indicate an active mobile data connection and its type.
1896
       ## Prerequisites for Mobile Data
1898
       For mobile data to work, the user must first have **cellular service**.
1899
       Refer to the "Understanding and Troubleshooting Your Phone's Cellular
       Service" guide if the user does not have service.
1900
1901
       ## Common Mobile Data Issues and Causes
1902
       Even with cellular service, mobile data problems might occur. Common
1903
       reasons include:
1904
           **Airplane Mode is ON**: Disables all wireless connections,
1905
       including mobile data.
1906
           **Mobile Data is Turned OFF**: The main switch for mobile data might
1907
       be disabled in the phone's settings.
1908
           **Roaming Issues (When User is Abroad) **:
               Data Roaming is turned OFF on the phone.
               The line is not roaming enabled.
1910
           **Data Plan Limits Reached**: The user may have used up their
1911
       monthly data allowance, and the carrier has slowed down or cut off data.
1912
           **Data Saver Mode is ON**: This feature restricts background data
1913
       usage and can make some apps or services seem slow or unresponsive to
1914
       save data.
           **VPN Issues**: An active VPN connection might be slow or
1915
       misconfigured, affecting data speeds or connectivity.
1916
          **Bad Network Preferences**: The phone is set to an older network
1917
       technology like 2G/3G.
1918
1919
       ## Diagnosing Mobile Data Issues
       `run_speed_test()` can be used to check for potential issues with mobile
1920
       data.
1921
       When mobile data is unavailable a speed test should return 'no
1922
       connection'.
1923
       If data is available, a speed test will also return the data speed.
1924
       Any speed below 'Excellent' is considered slow.
1925
       ## Troubleshooting Mobile Data Problems
1926
       ### Airplane Mode
1927
       Refer to the "Understanding and Troubleshooting Your Phone's Cellular
1928
       Service" section for instructions on how to check and turn off Airplane
1929
       Mode.
1930
       ### Mobile Data Disabled
1931
       Mobile data switch allows the phone to connect to the internet using the
1932
       carrier's cellular network.
1933
       If `check_network_status()` shows mobile data is disabled, guide the user to use `toggle_data()` to turn mobile data ON.
1934
1935
       ### Addressing Data Roaming Problems
1936
       Data roaming allows the user to use their phone's data connection in
1937
       areas outside their home network (e.g. when traveling abroad).
1938
       If the user is outside their carrier's primary coverage area (roaming)
```

and mobile data isn't working, guide them to use `toggle_roaming()` to

You should check that the line associated with the phone number the user

provided is roaming enabled. If it is not, the user will not be able to

use their phone's data connection in areas outside their home network.

Refer to the general policy for guidelines on enabling roaming.

ensure Data Roaming is ON.

1940

1941

1942

```
1945
       ### Data Saver Mode
       Data Saver mode is a feature that restricts background data usage and
1947
       can affect data speeds.
1948
       If `check_data_restriction_status()` shows "Data Saver mode is ON",
       guide the user to use `toggle_data_saver_mode()` to turn it OFF.
1949
1950
       ### VPN Connection Issues
1951
      VPN (Virtual Private Network) is a feature that encrypts internet
1952
       traffic and can help improve data speeds and security.
1953
       However in some cases, a VPN can cause speed to drop significantly.
       If `check_vpn_status()` shows "VPN is ON and connected" and performance
1954
       level is "Poor", guide the user to use `disconnect_vpn()` to disconnect
1955
       the VPN.
1956
1957
       ### Data Plan Limits Reached
1958
       Each plan specify the maxium data usage per month.
       If the user's data usage for a line associated with the phone number the
1959
       user provided exceeds the plan's data limit, data connectivity will be
1960
      lost.
1961
      The user has 2 options:
1962
       - Change to a plan with more data.
       - Add more data to the line by "refueling" data at a price per GB
       specified by the plan.
1964
       Refer to the general policy for guidelines on those options.
1965
1966
       ### Optimizing Network Mode Preferences
1967
      Network mode preferences are the settings that determine the type of
1968
      cellular network the phone will connect to.
      Using older modes like 2G/3G can significantly limit speed.
1969
       If `check_network_mode_preference()` shows "2G" or "3G", guide the user
1970
       to use `set_network_mode_preference(mode: str)` with the mode
1971
       "4g_5g_preferred" to allow the phone to connect to 5G.
1972
1973
       # Understanding and Troubleshooting MMS (Picture/Video Messaging)
       This section explains for agents how to troubleshoot Multimedia
1974
       Messaging Service (MMS), which allows users to send and receive messages
1975
       containing pictures, videos, or audio.
1976
1977
       ## What is MMS?
1978
       MMS is an extension of SMS (text messaging) that allows for multimedia
       content. When a user sends a photo to a friend via their messaging app,
1979
       they're typically using MMS.
1980
1981
       ## Prerequisites for MMS
1982
       For MMS to work, the user must have cellular service and mobile data
1983
       (any speed).
       Refer to the "Understanding and Troubleshooting Your Phone's Cellular
1984
       Service" and "Understanding and Troubleshooting Your Phone's Mobile
1985
       Data" sections for more information.
1986
1987
       ## Common MMS Issues and Causes
1988
          **No Cellular Service or Mobile Data Off/Not Working**: The most
       common reasons. MMS relies on these.
1989
          **Incorrect APN Settings**: Specifically, a missing or incorrect
1990
1991
          **Connected to 2G Network**: 2G networks are generally not suitable
1992
       for MMS.
          **Wi-Fi Calling Configuration**: In some cases, how Wi-Fi Calling is
       configured can affect MMS, especially if your carrier doesn't support
1994
       MMS over Wi-Fi.
1995
           **App Permissions**: The messaging app needs permission to access
```

storage (for the media files) and usually SMS functionalities.

1996

Diagnosing MMS Issues

'can_send_mms()' tool on the user's phone can be used to check if the user is facing an MMS issue.

Troubleshooting MMS Problems

Ensuring Basic Connectivity for MMS

Successful MMS messaging relies on fundamental service and data connectivity. This section covers verifying these prerequisites. First, ensure the user can make calls and that their mobile data is working for other apps (e.g., browsing the web). Refer to the "Understanding and Troubleshooting Your Phone's Cellular Service" and "Understanding and Troubleshooting Your Phone's Mobile Data" sections if needed.

Unsuitable Network Technology for MMS

MMS has specific network requirements; older technologies like 2G are insufficient. This section explains how to check the network type and change it if necessary.

MMS requires at least a 3G network connection; 2G networks are generally not suitable.

If `check_network_status()` shows "2G", guide the user to use `set_network_mode_preference(mode: str)` to switch to a network mode that includes 3G, 4G, or 5G (e.g., `"4g_5g_preferred"` or `"4g_only"`).

Verifying APN (MMSC URL) for MMS

MMSC is the Multimedia Messaging Service Center. It is the server that handles MMS messages. Without a correct MMSC URL, the user will not be able to send or receive MMS messages.

Those are specified as part of the APN settings. Incorrect MMSC URL, are a very common cause of MMS issues.

If `check_apn_settings()` shows MMSC URL is not set, guide the user to use `reset_apn_settings()` to reset the APN settings.

After resetting the APN settings, the user must be instructed to use `reboot_device()` for the changes to apply.

Investigating Wi-Fi Calling Interference with MMS

Wi-Fi Calling settings can sometimes conflict with MMS functionality. If `check_wifi_calling_status()` shows "Wi-Fi Calling is ON", guide the user to use `toggle_wifi_calling()` to turn it OFF.

Messaging App Lacks Necessary Permissions

The messaging app needs specific permissions to handle media and send messages.

If `check_app_permissions(app_name="messaging")` shows "storage" and "sms" permissions are not listed as granted, guide the user to use `grant_app_permission(app_name="messaging", permission="storage")` and `grant_app_permission(app_name="messaging", permission="sms")` to grant the necessary permissions.

In No-User mode, the agent is provided with a version of those policies that have been rephrased when needed. (e.g instructions like "ask user to do X", are rephrased as "perform action X")

D.2.3 TECHNICAL SUPPORT POLICY (WORKFLOW)

Introduction

This document serves as a comprehensive guide for technical support agents. It provides detailed procedures and troubleshooting steps to assist users experiencing common issues with their phone's cellular service, mobile data connectivity, and Multimedia Messaging Service (MMS). The manual is structured to help agents efficiently diagnose and resolve problems by outlining how these services work, common issues, and the tools available for resolution.

2052 2053 The main sections covered are: 2054 **Understanding and Troubleshooting Your Phone's Cellular Service**: 2055 Addresses issues related to network connection, signal strength, and SIM 2056 card problems. **Understanding and Troubleshooting Your Phone's Mobile Data**: 2057 Focuses on problems with internet access via the cellular network, 2058 including speed and connectivity. 2059 **Understanding and Troubleshooting MMS (Picture/Video Messaging) **: 2060 Covers issues related to sending and receiving multimedia messages. 2061 Make sure you try all the possible ways to resolve the user's issue 2062 before transferring to a human agent. 2063 2064 # What the user can do on their device 2065 Here are the actions a user is able to take on their device. 2066 You must understand those well since as part of technical support you will have to help the customer perform series of actions 2067 2068 ## Diagnostic Actions (Read-only) 2069 1. **check_status_bar** - Shows what icons are currently visible in your 2070 phone's status bar (the area at the top of the screen). 2071 - Airplane mode status ("Airplane Mode" when enabled) 2072 - Network signal strength ("No Signal", "Poor", "Fair", "Good", "Excellent") 2073 - Network technology (e.g., "5G", "4G", etc.) 2074 - Mobile data status ("Data Enabled" or "Data Disabled") 2075 - Data saver status ("Data Saver" when enabled) 2076 - Wi-Fi status ("Connected to [SSID]" or "Enabled") - VPN status ("VPN Connected" when connected) 2077 - Battery level ("[percentage]%") 2078 2. **check_network_status** - Checks your phone's connection status to cellular networks and Wi-Fi. Shows airplane mode status, signal 2079 2080 strength, network type, whether mobile data is enabled, and whether data 2081 roaming is enabled. Signal strength can be "none", "poor" (1bar), "fair" (2 bars), "good" (3 bars), "excellent" (4+ bars). 2082 3. **check_network_mode_preference** - Checks your phone's network mode 2083 preference. Shows the type of cellular network your phone prefers to 2084 connect to (e.g., 5G, 4G, 3G, 2G). 2085 4. **check_sim_status** - Checks if your SIM card is working correctly 2086 and displays its current status. Shows if the SIM is active, missing, or locked with a PIN or PUK code. 2087 5. **check_data_restriction_status** - Checks if your phone has any data-limiting features active. Shows if Data Saver mode is on and 2089 whether background data usage is restricted globally. 2090 6. **check_apn_settings** - Checks the technical APN settings your phone 2091 uses to connect to your carrier's mobile data network. Shows current APN name and MMSC URL for picture messaging. 2092 7. **check_wifi_status** - Checks your Wi-Fi connection status. Shows if 2093 Wi-Fi is turned on, which network you're connected to (if any), and the 2094 signal strength. 2095 8. **check_wifi_calling_status** - Checks if Wi-Fi Calling is enabled on 2096 your device. This feature allows you to make and receive calls over a Wi-Fi network instead of using the cellular network. 2097 9. **check vpn status** - Checks if you're using a VPN (Virtual Private 2098 Network) connection. Shows if a VPN is active, connected, and displays 2099 any available connection details. 2100 10. **check_installed_apps** - Returns the name of all installed apps on 2101 the phone. 11. **check_app_status** - Checks detailed information about a specific 2102 app. Shows its permissions and background data usage settings. 2103 12. **check_app_permissions** - Checks what permissions a specific app 2104 currently has. Shows if the app has access to features like storage,

2105

camera, location, etc.

```
13. **run_speed_test** - Measures your current internet connection speed
2107
       (download speed). Provides information about connection quality and what
2108
       activities it can support. Download speed can be "unknown", "very poor",
2109
       "poor", "fair", "good", or "excellent".
2110
       14. **can_send_mms** - Checks if the messaging app can send MMS
       messages.
2111
2112
       ## Fix Actions (Write/Modify)
2113
       1. **set_network_mode_preference** - Changes the type of cellular
2114
       network your phone prefers to connect to (e.g., 5G, 4G, 3G).
2115
       Higher-speed networks (5G, 4G) provide faster data but may use more
2116
       2. **toggle_airplane_mode** - Turns Airplane Mode ON or OFF. When ON, it
2117
       disconnects all wireless communications including cellular, Wi-Fi, and
2118
       Bluetooth.
2119
       3. **reseat_sim_card** - Simulates removing and reinserting your SIM
2120
       card. This can help resolve recognition issues.
       4. **toggle_data** - Turns your phone's mobile data connection ON or
2121
       OFF. Controls whether your phone can use cellular data for internet
2122
       access when Wi-Fi is unavailable.
2123
       5. **toggle_roaming** - Turns Data Roaming ON or OFF. When ON, roaming
2124
       is enabled and your phone can use data networks in areas outside your
2125
       carrier's coverage.
       6. **toggle_data_saver_mode** - Turns Data Saver mode ON or OFF. When
2126
       ON, it reduces data usage, which may affect data speed.
2127
       7. **set_apn_settings** - Sets the APN settings for the phone.
2128
       8. **reset_apn_settings** - Resets your APN settings to the default
2129
       settings.
2130
       9. **toggle_wifi** - Turns your phone's Wi-Fi radio ON or OFF. Controls
       whether your phone can discover and connect to wireless networks for
2131
       internet access.
2132
       10. **toggle_wifi_calling** - Turns Wi-Fi Calling ON or OFF. This
2133
       feature allows you to make and receive calls over Wi-Fi instead of the
2134
       cellular network, which can help in areas with weak cellular signal.
2135
       11. **connect_vpn** - Connects to your VPN (Virtual Private Network).
       12. **disconnect_vpn** - Disconnects any active VPN (Virtual Private
2136
       Network) connection. Stops routing your internet traffic through a VPN
2137
       server, which might affect connection speed or access to content.
2138
       13. **grant_app_permission** - Gives a specific permission to an app
2139
       (like access to storage, camera, or location). Required for some app
2140
       functions to work properly.
       14. **reboot_device** - Restarts your phone completely. This can help
2141
       resolve many temporary software glitches by refreshing all running
2142
       services and connections.
2143
2144
       # Understanding and Troubleshooting Your Phone's Cellular Service
2145
       This section details for agents how a user's phone connects to the
       cellular network (often referred to as "service") and provides
2146
       procedures to troubleshoot common issues. Good cellular service is
2147
       required for calls, texts, and mobile data.
2148
2149
       ## Common Service Issues and Their Causes
2150
       If the user is experiencing service problems, here are some common
       causes:
2151
2152
           **Airplane Mode is ON**: This disables all wireless radios,
2153
       including cellular.
2154
           **SIM Card Problems**:
2155
               Not inserted or improperly seated.
               Locked due to incorrect PIN/PUK entries.
2156
           **Incorrect Network Settings**: APN settings might be incorrect
2157
       resulting in a loss of service.
2158
```

Carrier Issues: Your line might be inactive due to billing

2159

problems.

2160 2161 2162 ## Diagnosing Service Issues 2163 `check_status_bar()` can be used to check if the user is facing a 2164 service issue. If there is cellular service, the status bar will return a signal 2165 strength indicator. 2166 2167 ## Troubleshooting Service Problems 2168 ### Airplane Mode 2169 Airplane Mode is a feature that disables all wireless radios, including cellular. If it is enabled, it will prevent any cellular connection. 2170 You can check if Airplane Mode is ON by using `check_status_bar()` or 2171 `check_network_status()`. 2172 If it is ON, guide the user to use `toggle_airplane_mode()` to turn it 2173 2174 ### SIM Card Issues 2175 The SIM card is the physical card that contains the user's information 2176 and allows the phone to connect to the cellular network. Problems with the SIM card can lead to a complete loss of service. 2178 The most common issue is that the SIM card is not properly seated or the 2179 user has entered the wrong PIN or PUK code. Use `check_sim_status()` to check the status of the SIM card. 2180 If it shows "Missing", guide the user to use `reseat_sim_card()` to ensure the SIM card is correctly inserted. 2182 If it shows "Locked" (due to incorrect PIN or PUK entries), **escalate 2183 to technical support for assistance with SIM security**. 2184 If it shows "Active", the SIM itself is likely okay. 2185 ### Incorrect APN Settings 2186 Access Point Name (APN) settings are crucial for network connectivity. 2187 If `check_apn_settings()` shows "Incorrect", guide the user to use 2188 `reset_apn_settings()` to reset the APN settings. 2189 After resetting the APN settings, the user must be instructed to use `reboot_device()` for the changes to apply. 2190 2191 ### Line Suspension 2192 If the line is suspended, the user will not have cellular service. 2193 Investigate if the line is suspended. Refer to the general agent policy 2194 for guidelines on handling line suspensions. If the line is suspended and the agent can lift the suspension (per 2195 general policy), verify if service is restored. 2196 If the suspension cannot be lifted by the agent (e.g., due to 2197 contract end date as mentioned in general policy, or other reasons not 2198 resolvable by the agent), **escalate to technical support**. 2199 2200 # Understanding and Troubleshooting Your Phone's Mobile Data 2201 This section explains for agents how a user's phone uses mobile data for 2202 internet access when Wi-Fi is unavailable, and details troubleshooting 2203 for common connectivity and speed issues. 2204 2205 ## What is Mobile Data? Mobile data allows the phone to connect to the internet using the 2206 carrier's cellular network. This enables browsing websites, using apps, 2207 streaming video, and sending/receiving emails when not connected to 2208 Wi-Fi. The status bar usually shows icons like "5G", "LTE", "4G", "3G", 2209 "H+", or "E" to indicate an active mobile data connection and its type. 2210 ## Prerequisites for Mobile Data 2211 For mobile data to work, the user must first have **cellular service**.

Refer to the "Understanding and Troubleshooting Your Phone's Cellular

Service" guide if the user does not have service.

2212

```
2214
2215
       ## Common Mobile Data Issues and Causes
       Even with cellular service, mobile data problems might occur. Common
2217
       reasons include:
2218
           **Airplane Mode is ON**: Disables all wireless connections,
2219
       including mobile data.
2220
           **Mobile Data is Turned OFF**: The main switch for mobile data might
2221
       be disabled in the phone's settings.
2222
           **Roaming Issues (When User is Abroad) **:
2223
               Data Roaming is turned OFF on the phone.
               The line is not roaming enabled.
2224
           **Data Plan Limits Reached**: The user may have used up their
2225
       monthly data allowance, and the carrier has slowed down or cut off data.
2226
          **Data Saver Mode is ON**: This feature restricts background data
2227
       usage and can make some apps or services seem slow or unresponsive to
2228
       save data.
          **VPN Issues**: An active VPN connection might be slow or
2229
       misconfigured, affecting data speeds or connectivity.
2230
          **Bad Network Preferences**: The phone is set to an older network
2231
       technology like 2G/3G.
2232
2233
       ## Diagnosing Mobile Data Issues
       run_speed_test() can be used to check for potential issues with mobile
2234
2235
       When mobile data is unavailable a speed test should return 'no
2236
       connection'.
2237
       If data is available, a speed test will also return the data speed.
2238
       Any speed below 'Excellent' is considered slow.
2239
       ## Troubleshooting Mobile Data Problems
2240
       ### Airplane Mode
2241
       Refer to the "Understanding and Troubleshooting Your Phone's Cellular
2242
       Service" section for instructions on how to check and turn off Airplane
2243
       Mode.
2244
       ### Mobile Data Disabled
2245
       Mobile data switch allows the phone to connect to the internet using the
2246
       carrier's cellular network.
2247
       If `check_network_status()` shows mobile data is disabled, guide the
2248
       user to use `toggle_data()` to turn mobile data ON.
2249
       ### Addressing Data Roaming Problems
2250
       Data roaming allows the user to use their phone's data connection in
2251
       areas outside their home network (e.g. when traveling abroad).
2252
      If the user is outside their carrier's primary coverage area (roaming)
2253
       and mobile data isn't working, quide them to use `toggle_roaming()` to
       ensure Data Roaming is ON.
       You should check that the line associated with the phone number the user
2255
       provided is roaming enabled. If it is not, the user will not be able to
2256
       use their phone's data connection in areas outside their home network.
2257
       Refer to the general policy for guidelines on enabling roaming.
2258
       ### Data Saver Mode
2259
       Data Saver mode is a feature that restricts background data usage and
2260
       can affect data speeds.
2261
       If `check_data_restriction_status()` shows "Data Saver mode is ON",
2262
       guide the user to use `toggle_data_saver_mode()` to turn it OFF.
2263
       ### VPN Connection Issues
2264
       VPN (Virtual Private Network) is a feature that encrypts internet
2265
       traffic and can help improve data speeds and security.
2266
       However in some cases, a VPN can cause speed to drop significantly.
2267
       If `check_vpn_status()` shows "VPN is ON and connected" and performance
```

level is "Poor", guide the user to use `disconnect_vpn()` to disconnect

the VPN.

2268 2269 ### Data Plan Limits Reached Each plan specify the maxium data usage per month. 2271 If the user's data usage for a line associated with the phone number the 2272 user provided exceeds the plan's data limit, data connectivity will be lost. 2273 The user has 2 options: 2274 - Change to a plan with more data. 2275 - Add more data to the line by "refueling" data at a price per GB 2276 specified by the plan. Refer to the general policy for guidelines on those options. 2277 2278 ### Optimizing Network Mode Preferences 2279 Network mode preferences are the settings that determine the type of 2280 cellular network the phone will connect to. 2281 Using older modes like 2G/3G can significantly limit speed. 2282 If `check_network_mode_preference()` shows "2G" or "3G", guide the user to use `set_network_mode_preference(mode: str)` with the mode 2283 "4g_5g_preferred" to allow the phone to connect to 5G. 2284 2285 # Understanding and Troubleshooting MMS (Picture/Video Messaging) 2286 This section explains for agents how to troubleshoot Multimedia 2287 Messaging Service (MMS), which allows users to send and receive messages 2288 containing pictures, videos, or audio. 2289 ## What is MMS? 2290 ${
m MMS}$ is an extension of ${
m SMS}$ (text messaging) that allows for multimedia 2291 content. When a user sends a photo to a friend via their messaging app, they're typically using MMS. 2293 ## Prerequisites for MMS 2294 For MMS to work, the user must have cellular service and mobile data 2295 (any speed). 2296 Refer to the "Understanding and Troubleshooting Your Phone's Cellular 2297 Service" and "Understanding and Troubleshooting Your Phone's Mobile Data" sections for more information. 2298 2299 ## Common MMS Issues and Causes 2300 **No Cellular Service or Mobile Data Off/Not Working**: The most 2301 common reasons. MMS relies on these. 2302 **Incorrect APN Settings**: Specifically, a missing or incorrect MMSC URL. 2303 **Connected to 2G Network**: 2G networks are generally not suitable 2304 for MMS. 2305 **Wi-Fi Calling Configuration**: In some cases, how Wi-Fi Calling is 2306 configured can affect MMS, especially if your carrier doesn't support 2307 MMS over Wi-Fi. **App Permissions**: The messaging app needs permission to access 2308 storage (for the media files) and usually SMS functionalities. 2309 2310 ## Diagnosing MMS Issues 2311 can send mms() tool on the user's phone can be used to check if the 2312 user is facing an MMS issue. 2313 ## Troubleshooting MMS Problems 2314 ### Ensuring Basic Connectivity for MMS 2315 Successful MMS messaging relies on fundamental service and data 2316 connectivity. This section covers verifying these prerequisites. 2317 First, ensure the user can make calls and that their mobile data is working for other apps (e.g., browsing the web). Refer to the 2318 "Understanding and Troubleshooting Your Phone's Cellular Service" and

"Understanding and Troubleshooting Your Phone's Mobile Data" sections if

2319

2320

2321

needed.

```
2322
       ### Unsuitable Network Technology for MMS
2323
       MMS has specific network requirements; older technologies like 2G are
2324
       insufficient. This section explains how to check the network type and
2325
       change it if necessary.
2326
       MMS requires at least a 3G network connection; 2G networks are generally
       not suitable.
2327
       If `check_network_status()` shows "2G", guide the user to use
   `set_network_mode_preference(mode: str)` to switch to a network mode
2328
2329
       that includes 3G, 4G, or 5G (e.g., "4g_5g_preferred" or "4g_only").
2330
2331
       ### Verifying APN (MMSC URL) for MMS
       MMSC is the Multimedia Messaging Service Center. It is the server that
2332
       handles MMS messages. Without a correct MMSC URL, the user will not be
2333
       able to send or receive MMS messages.
2334
       Those are specified as part of the APN settings. Incorrect MMSC URL, are
2335
       a very common cause of MMS issues.
2336
       If `check_apn_settings()` shows MMSC URL is not set, guide the user to
       use `reset_apn_settings()` to reset the APN settings.
2337
       After resetting the APN settings, the user must be instructed to use
2338
       `reboot_device()` for the changes to apply.
2339
2340
       ### Investigating Wi-Fi Calling Interference with MMS
       Wi-Fi Calling settings can sometimes conflict with MMS functionality.
       If `check_wifi_calling_status()` shows "Wi-Fi Calling is ON", guide the
2342
       user to use `toggle_wifi_calling()` to turn it OFF.
2343
2344
       ### Messaging App Lacks Necessary Permissions
2345
       The messaging app needs specific permissions to handle media and send
2346
       messages.
       If `check_app_permissions(app_name="messaging")` shows "storage" and
2347
       "sms" permissions are not listed as granted, quide the user to use
2348
       `grant_app_permission(app_name="messaging", permission="storage")` and
2349
        `grant_app_permission(app_name="messaging", permission="sms")` to grant
2350
       the necessary permissions.
2351
```

D.2.4 TROUBLESHOOTING WORKFLOW GRAPHS

To help the agent understand the troubleshooting workflow, we provide a decision graph for each issue type.

E USER SIMULATOR QUALITY

2352 2353

23542355

2356

2357

23592360

23612362

2363

2364

2365

2366

2367

2368 2369

2370

2371

2372

2373

2374

2375

E.1 COMMON ERROR TYPES AND FAILURE MODES (RETAIL)

Manual analysis of the 20 annotated errors in the retail domain exposes three recurring failure modes:

- Conversation-structure rule violation (11/20) the simulator breaks turn-taking or dialogue-flow instructions (e.g., mixes tool calls with natural language in the same turn).
- **Premature termination** (3/20) the simulator halts the conversation immediately after the user's confirmation (###STOP###), preventing the agent from completing the transaction.
- **Ungrounded reference** (2/20) the simulator invents or misstates contextual details such as payment method or order status.
- **Missing constraint** (4/20) the simulator omits a required instruction (e.g., neglects to request an alternative SKU when the desired colour is unavailable).

Most task-critical errors stem from either premature termination or missing constraints, whereas conversation-structure violations and ungrounded references are typically task-benign and readily recoverable by the agent.

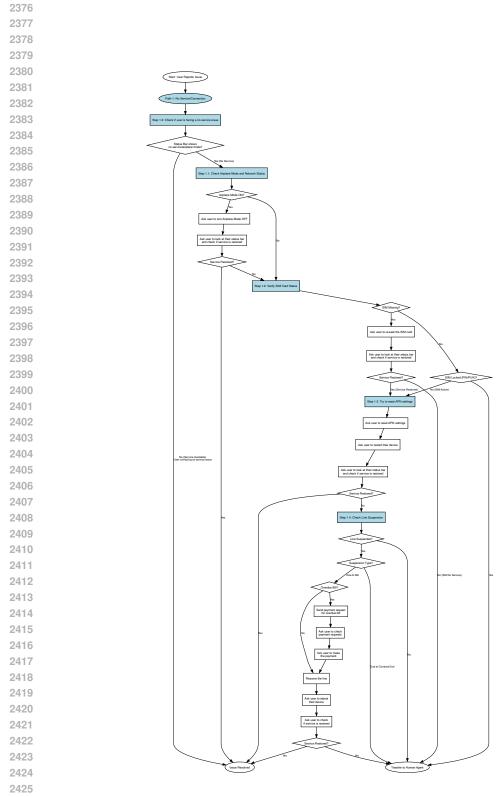


Figure 9: Troubleshooting workflow for service_issue

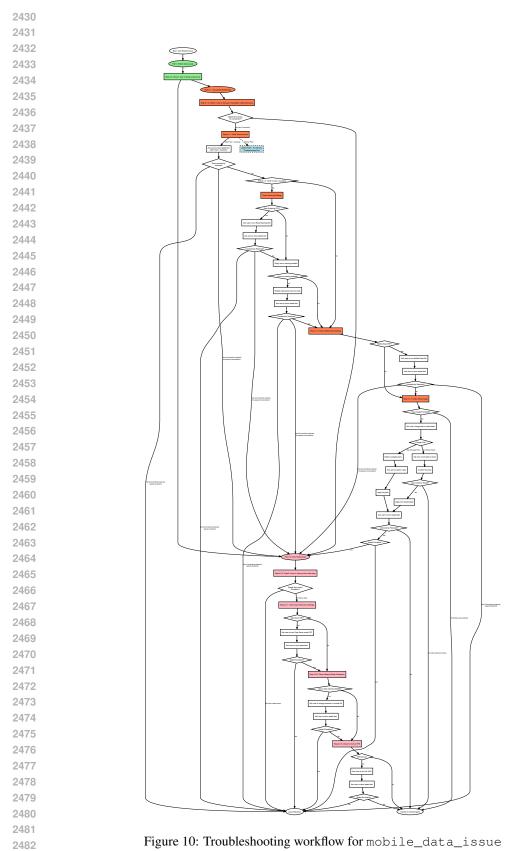


Figure 10: Troubleshooting workflow for mobile_data_issue

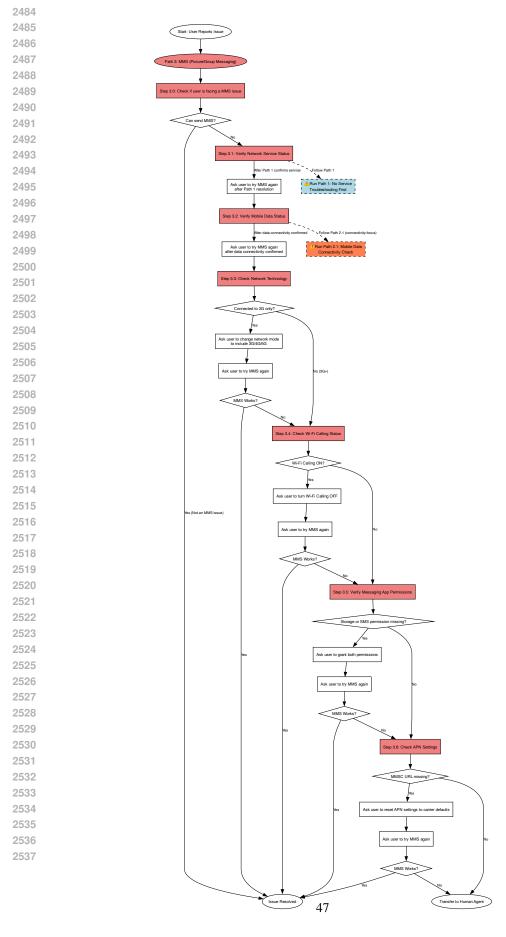


Figure 11: Troubleshooting workflow for mms_issue

E.2 COMMON ERROR TYPES AND FAILURE MODES (AIRLINE) Preliminary analysis of annotated airline dialogues identifies three dominant issues: • Missing constraint (11/47)

- Conversation-structure rule violation (19/47)
 Ungrounded reference (15/47)
- oligiounded reference (13/47)

• Premature termination (2/47)

Again, most task-critical errors stem from missing constraints.

E.3 COMMON ERROR TYPES AND FAILURE MODES (TELECOM)

Manual analysis of the 8 annotated errors in the telecom domain reveals that all errors are caused by **Premature termination** (8/8, 100%). The user returns the ###TRANSFER### token before the assistant calls the transfer tool. This error is critical only if this causes a task failure by preventing the agent from performing an expected transfer.

E.4 ILLUSTRATIVE CASES

Task-critical error (retail). Premature termination. After the assistant requests final confirmation to modify order #W5061109, the simulator ends the dialogue with ###STOP### immediately following confirming its intent, preventing the agent from executing the modification and issuing a receipt:

assistant: Please confirm that you have provided all the items you want to modify in this order, as this action can only be done once...

user: Yes, I confirm that I want to modify my pending order #W5061109 to ...Please proceed. Thank you!###STOP###

Task-critical error (airline). Missing constraint. Here the simulator is instructed to ask for a transfer to a human agent in the case of outbound trip unavailability, but proceeds to book the return trip instead, violating goal constraints.

assistant: It appears that one of your original outbound flights (HAT030 on 2024-05-13) is no longer available in economy class for the upgrade. To proceed, I will need to find available economy class flights for your outbound journey as well. Would you like me to search for available economy class flights for your outbound trip (Los Angeles \rightarrow Atlanta) and your new return trip (Atlanta \rightarrow Las Vegas)? Please confirm so I can continue.

user: Yes, please go ahead and find available economy class flights for both the outbound trip to Atlanta and the new return trip to Las Vegas. Thank you! ###STOP###

Task-critical error (telecom). Premature termination. The simulator ends the dialogue with ###TRANSFER### before the agent could perform the actual transfer:

assistant: ... The next step is to escalate your issue to a human technical support agent who can provide more advanced help, possibly including SIM replacement or device diagnostics. Would you like me to transfer you to a human agent for further assistance?

user: Yes, please transfer me to someone who can help further. I just want my phone to work again. Thank you for your help so far.###TRANSFER###

F USE OF LARGE LANGUAGE MODELS

As part of the paper writing process, we utilized Large Language Models (LLMs). The use of these models was restricted to improving the clarity and readability of the text. Specifically, we used LLMs

for grammar correction, rephrasing sentences for better flow, and ensuring consistent terminology. The core ideas, experimental results, and scientific contributions presented in this paper are entirely our own.