

---

# UserBench: An Interactive Gym Environment for User-Centric Agents

---

Cheng Qian<sup>1,2</sup>, Zuxin Liu<sup>1</sup>, Akshara Prabhakar<sup>1</sup>, Zhiwei Liu<sup>1</sup>, Jianguo Zhang<sup>1</sup>,  
Haolin Chen<sup>1</sup>, Heng Ji<sup>2</sup>, Weiran Yao<sup>1</sup>, Shelby Heinecke<sup>1</sup>, Silvio Savarese<sup>1</sup>,  
Caiming Xiong<sup>1</sup>, Huan Wang<sup>1</sup>

<sup>1</sup>Salesforce AI Research <sup>2</sup>University of Illinois Urbana-Champaign

## Abstract

Large Language Models (LLMs)-based agents have made impressive progress in reasoning and tool use, enabling them to solve complex tasks. However, their ability to proactively collaborate with users, especially when goals are vague, evolving, or indirectly expressed, remains underexplored. To address this gap, we introduce **UserBench**, a user-centric benchmark designed to evaluate agents in multi-turn, preference-driven interactions. UserBench features simulated users who start with underspecified goals and reveal preferences incrementally, requiring agents to proactively clarify intent and make grounded decisions with tools. Our evaluation of leading open- and closed-source LLMs reveals a significant disconnect between task completion and user alignment. For instance, models provide answers that fully align with all user intents only 20% of the time on average, and even the most advanced models uncover fewer than 30% of all user preferences through active interaction. These results highlight the challenges of building agents that are not just capable task executors, but true collaborative partners. UserBench offers an interactive environment to measure and advance this critical capability.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in complex reasoning [12, 9], code generation [11, 8, 39], and solving advanced mathematical problems [35, 30, 29]. Their limitations in updated knowledge and precise computation are increasingly mitigated through tool use [18, 25, 38, 21] and tool creation [22, 2, 37, 26], enabling them to retrieve information and interact with external environments including database, web and games [5, 23, 42]. This tool-augmented reasoning grants LLMs agentic capabilities: they can autonomously execute tasks via these tools.

However, existing agentic environments often overlook a critical dimension: the role of the user. Despite achieving strong task performance, agents frequently fail to satisfy real user needs due to their inability to understand, adapt to, and collaborate with the task initiator [24, 20]. Current evaluations primarily assess tool use and task execution, rarely considering whether the agent effectively interprets and aligns with the user’s underlying and evolving intent.

This motivates our central research question: **How can we evaluate agents from a user-centric perspective?** To answer this, we first examine how users typically communicate goals. Human communication is inherently a joint activity, where meaning is co-constructed through interaction [4]. Moreover, language is inherently ambiguous, making it difficult for users to fully and clearly convey their intent in a single interaction. [17]. As such, user instructions tend to share three core traits: (i) **Underspecification**: users often initiate requests before fully formulating their goals; (ii)

Benchmark	Multi-turn Interaction	Goal Ambiguity	Tool Use	Dynamic State	Multi-Aspect Reasoning	User Simulation	Domain Diversity	Custom- izable	Scal- able
$\tau^2$ -Bench[1]	✓	✗	✓	✓	✗	✓	✗	✓	✓
$\tau$ -Bench[36]	✓	✗	✓	✓	✓	✓	✓	✓	✓
ToolSandbox[19]	✓	✓	✓	✓	✓	✓	✓	✓	✓
MINT[34]	✓	✓	✓	✓	✓	✓	✓	✓	✓
IN3[24]	✓	✓	✗	✗	✗	✓	✓	✗	✗
ToolTalk[6]	✓	✗	✓	✗	✗	✗	✓	✗	✗
API-Bank[14]	✗	✗	✓	✗	✗	✗	✓	✓	✓
MetaTool[10]	✗	✗	✓	✗	✗	✗	✓	✓	✓
<b>UserBench (Ours)</b>	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 1: For each existing benchmark, the table indicates whether the corresponding trait is fully addressed (✓), partially addressed (✓), or not addressed (✗). We provide detailed explanation for each trait being compared in Appendix A.

**Incrementality:** intent emerges and evolves across interaction; and (iii) **Indirectness:** users may obscure or soften their true intent due to social or strategic reasons.

Revolving around these traits, we introduce **UserBench**, a user-centric environment designed to facilitate an agent’s ability to engage in meaningful, multi-turn interactions with users who exhibit these traits. In UserBench, simulated users provide initial vague task instruction (underspecification), gradually reveal preferences over time (incrementality), and often do so implicitly (indirectness). Agents must proactively clarify goals, interpret subtle cues, and adaptively reason through tools.

Built on the standard Gymnasium framework, UserBench offers a modular, extensible setup with a standardized interaction interface and a stable tool-use backend, enabling rigorous and reproducible evaluation. We benchmark several leading open- and closed-source models and find that current LLMs still struggle to interactively uncover and act on user preferences. For instance, scores drop by over 40% on average when models are restricted to selecting only one option per traveling aspect in UserBench, revealing their difficulty in making optimal decisions. Moreover, models provide answers that fully align with all user intents only 20% of the time on average, and even the best-performing models elicit less than 30% of all user preferences through active querying, suggesting limited ability to engage in purposeful, user-driven dialogue. While strong models handle tool use reliably, they remain brittle in understanding implicit and nuanced human needs, highlighting the core challenge that UserBench aims to measure. We summarize our contributions as follows:

- **Data:** We propose 4K+ scenarios that capture grounded communication challenges (underspecification, incrementality, indirectness) through a carefully curated pipeline.
- **Environment:** We introduce UserBench, a scalable and modular gym environment designed both as a benchmark and a training ground for LLM agents engaged in multi-turn, preference-driven user interactions.
- **Analysis:** We show that despite strong tool use capability, current models still struggle to comprehensively uncover and align with user intent, revealing limits in user-centric reasoning.

We view this work as a foundational step toward truly user-centric agents: not just efficient executors, but collaborative teammates capable of aligning with nuanced human intent.

## 2 Related Work

**User-centric environments for LLM evaluation.** Recent work has increasingly emphasized the importance of evaluating LLMs under realistic, user-centric conditions. Traditional benchmarks assess task success assuming fully specified prompts, overlooking how users often communicate underspecified, incremental, or indirect goals. To address this, benchmarks like User Reported Scenarios [33], Intention-in-Interaction [24], and WildBench [16] compile real user queries to evaluate whether models align with nuanced human intent. These datasets emphasize preference satisfaction, multi-intent understanding, and performance in “in-the-wild” conditions. Complementing this, other benchmarks such as MINT [34], PrefEval [41],  $\tau$ -Bench [36], and  $\tau^2$ -Bench [1] focus on dynamic, multi-turn interactions, testing whether agents can incorporate feedback, handle evolving preferences, and maintain user alignment over time. However, several limitations remain. For example, a large portion of interactions are automatically synthesized in these benches, often resulting in overly lengthy and unnatural goal formulations that diverge from how users typically express themselves. In contrast, UserBench explicitly models the three core interactive traits of user communication that are critical for real-world alignment. It also spans hundreds of diverse goal configurations, enabling fine-grained

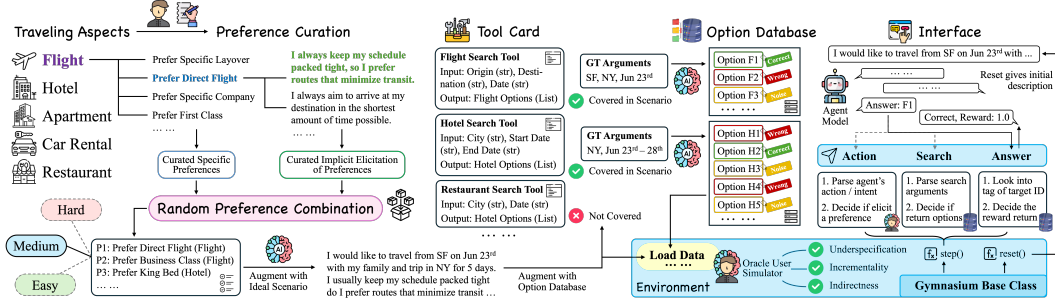


Figure 1: The pipeline of UserBench Construction, including data gathering, preference curation, tool augmentation, environment setup and interface design. *Random Preference Combinations* makes our pipeline highly scalable and enables the efficient generation of even 10K+ data points.

evaluation on a wide range of user intents, thus providing a more comprehensive and behaviorally grounded setting for assessing user-centric agents. We provide a comprehensive comparison of our gym environment with recent environments and benchmarks in Table 1.

**User-centric agent designs.** Designing agents that collaborate effectively with users requires modeling ambiguity, evolving intent, and user-specific preferences. Standard instruction-tuned models often hallucinate intent or avoid asking clarifying questions. Recent work aims to teach agents to proactively clarify underspecified instructions. For instance, models trained with simulated ambiguity resolution [40, 3] better recognize when to ask versus answer. Other work tackles user modeling directly: TravelPlanner+ [31] and PRELUDE [7] build agents that personalize responses using either explicit profiles or latent preferences inferred from user edits. While these designs make progress toward adaptive, user-aware agents, they often target narrow personalization tasks or rely on static user models; in contrast, our work provides a general-purpose evaluation environment that systematically tests an agent’s ability to dynamically uncover and respond to emergent user intent across varied and customizable interaction scenarios.

### 3 UserBench

We build UserBench upon the standardized gymnasium environment, focusing on travel planning tasks where users gradually reveal their preferences in implicit ways. Compared to prior work, UserBench features a significantly more diverse and grounded collection of user preferences, each paired with carefully curated implicit expressions, and supports tool-augmented search with injected noise and cost-awareness. This setting challenges agents to perform search and reasoning in a user-centric manner, under constraints of partial, evolving, and indirect preference information. We describe the gym design details in the following.

#### 3.1 Data Gathering

As shown in Figure 1, we begin by identifying five core aspects of travel planning: flight, hotel, apartment, car rental, and restaurant. Our data gathering is guided by three key principles: *realism*, *diversity*, and *implicitness*. First, we ensure realism by curating preferences that reflect plausible, real-world user needs—such as preferring direct flights (flight) or king-sized beds (hotel). Second, we aim for diversity by collecting around a hundred distinct preferences that span travel aspects. Third, we emphasize implicitness by pairing each preference with multiple naturally phrased statements that express the intent indirectly. For instance, the preference for direct flights might be expressed as: “I always keep my schedule packed tight, so I prefer travel routes that minimize transit time.” These implicit expressions are carefully curated and validated to ensure that (i) they do not reveal the preference explicitly, and (ii) they still clearly imply the intended preference without ambiguity.

These curated preferences and their associated implicit expressions serve as the foundation for generating travel scenarios. We randomly sample and combine preferences across the five aspects, categorizing the resulting data into three difficulty tiers based on the number of preferences involved. For each combination, we generate an ideal travel arrangement statement that incorporates all relevant preferences using GPT-4o. This finally forms a data point in our environment.

Difficulty Tier	# Test	# Train	Composition	Global Dataset Metadata		Travel Aspect	# Scenarios	Preferences
Easy	118	666	Travel-22, 2222	Travel Aspects	5	Hotel	187 (39.7%)	15
Medium	201	1130	Travel-33, 233, 333	Pref. Categories	24	Restaurant	298 (63.3%)	19
Hard	152	855	Travel-44, 334, 444	Total Prefs.	82	Flight	275 (58.4%)	16
Total	471	2651	UserBench Test / Train	Elicitation Ways	902	Rental Car	243 (51.6%)	14
				Options (B/C/W/N)	1/3/10/5	Apartment	172 (36.5%)	18

Table 2: **UserBench Statistics:** Difficulty-tiered statistics (left), global configuration metadata (middle), and travel aspect coverage (right). *Pref.* denotes preference; *B/C/W/N* stands for Best, Correct, Wrong, and Noise options respectively; *# Scenarios* indicates the number of data scenarios in which each travel aspect appears.

### 3.2 Tool Augmentation

For each travel planning aspect, we develop a corresponding tool that simulates a database search. For example, the flight search tool accepts origin, destination, and date, returning a list of mock options. Instead of querying real-time data, we rely on pre-generated options to ensure stable and controlled tool outputs. This design choice is justified because (i) it guarantees consistency and quality of returned results, and (ii) our focus is user-centric reasoning rather than challenges of searching.

Each data point is associated with a tailored database of options for all relevant travel planning aspects. Specifically, for each aspect, we generate three types of options: correct (satisfying all preferences), wrong (violating at least one preference), and noise (incomplete or irrelevant to the query, such as a flight option with a destination different from the one the user searched for). These options are mixed to form the complete search space, with each data point containing up to over 100 options. Option generation is handled by GPT-4o under human supervision to ensure quality and diversity.

### 3.3 Environment Building

Each data point in the environment contains a travel scenario, a set of implicit user preferences, and a database of relevant options. The environment simulates an oracle user with access to these internal states, while engaging in multi-turn interactions with a agent model under test. Upon reset, the environment provides only basic trip information and high-level requirements (e.g., hotel or flight), without disclosing any of its specific preferences. Preferences are then elicited over time, either (i) when the tested model explicitly asks about a relevant aspect, or (ii) after a fixed number of turns without progress (see Appendix B for details), at which point a preference is randomly and proactively revealed. All elicited preferences are presented implicitly, in line with the associated data.

This design captures the three key traits and challenges we previously identify: *underspecification*, where only coarse needs are initially provided; *incrementality*, where user preferences emerge progressively as the conversation evolves; and *indirectness*, where preferences are stated implicitly through natural, carefully-curated, heuristic-based language. These properties reflect real-world user behaviors and make UserBench a robust testbed for user-centric reasoning.

### 3.4 Agent Interaction

We provide a standardized interaction interface for agents, supporting three types of actions: action, search, and answer. Through search, agents issue queries specifying the travel planning aspect and arguments (e.g., hotel with specific dates), which are matched against ground truth. If the arguments match, the environment returns the mixed set of relevant options; otherwise, it returns nothing. Through action, the agent communicates with the user, typically by asking clarifying questions. The environment interprets the intent and responds naturally, potentially revealing a preference in an implicit form. Finally, through answer, the agent selects its recommended option ID(s), which are evaluated against the ground-truth correct options. Since ideal itineraries often span multiple aspects, correct answers may include several options.

### 3.5 Variants, Scalability and Statistics

UserBench supports several extensions. First, it can simulate a noisy search environment where query results are corrupted at a configurable frequency, mimicking real-world web or tool errors. Second, it allows for budget-aware decision-making by rewarding models that not only find correct options but

Model Name	Best Exist Rate	Correct Exist Rate	Score	Search Attempt Valid (%)	Action Attempt Valid (%)	Preference Elicited (%) (Active / Passive)
GPT-4o	0.204	<b>0.361</b>	<b>0.329</b>	82.48	27.82	27.32 (24.06 / 3.26)
Gemini-2.5-Pro	0.245	0.328	0.317	77.40	29.26	29.71 (23.85 / 5.85)
Claude-4-Sonnet	<b>0.260</b>	0.318	0.307	74.28	24.26	<b>34.25 (26.31 / 7.94)</b>
Deepseek-V3	0.148	0.218	0.210	61.22	23.62	31.82 (18.90 / 12.92)
Qwen3-14B	0.157	0.197	0.209	70.56	<b>33.71</b>	30.07 (15.95 / 14.12)
Qwen3-32B	0.154	0.211	0.206	79.42	26.89	22.01 (13.04 / 8.96)
Llama-3.3-70B	0.120	0.209	0.198	60.51	11.29	23.25 (13.05 / 10.20)
Qwen3-8B	0.140	0.186	0.180	68.94	23.20	28.86 (11.15 / 17.70)
GPT-4o-mini	0.102	0.173	0.166	76.10	13.55	25.43 (12.91 / 12.52)
Llama-3.1-8B	0.073	0.171	0.159	57.99	12.13	19.87 (8.09 / 11.78)
Gemini-2.5-Flash	0.047	0.080	0.125	<b>83.62</b>	28.14	14.00 (11.15 / 2.86)

Table 3: UserBench main evaluation results across different models (single-choice setting).

Model Name	Best Exist Rate	Correct Exist Rate	Score	Search Attempt Valid (%)	Action Attempt Valid (%)	Preference Elicited (%) (Active / Passive)
GPT-4o	<b>0.652</b>	<b>0.725</b>	<b>0.710</b>	85.96	37.89	15.10 (13.11 / 1.99)
Gemini-2.5-Pro	0.604	0.685	0.673	80.43	38.01	37.25 (28.80 / 8.46)
Claude-4-Sonnet	0.586	0.619	0.612	<b>87.94</b>	<b>41.17</b>	<b>42.21 (30.67 / 11.54)</b>
Qwen3-32B	0.343	0.416	0.411	82.76	34.24	17.72 (10.81 / 6.90)
Deepseek-V3	0.351	0.398	0.391	61.91	21.06	16.96 (10.18 / 6.78)
Llama-3.3-70B	0.294	0.366	0.372	63.46	21.72	31.72 (16.68 / 15.04)
Gemini-2.5-Flash	0.279	0.316	0.311	84.27	33.58	15.32 (11.67 / 3.64)
GPT-4o-mini	0.227	0.294	0.283	81.54	24.90	20.18 (10.46 / 9.72)
Qwen3-14B	0.204	0.249	0.258	71.48	35.06	26.21 (13.07 / 13.14)
Llama-3.1-8B	0.117	0.252	0.234	61.38	12.20	26.01 (7.80 / 18.20)
Qwen3-8B	0.147	0.186	0.182	69.83	23.26	20.31 (7.95 / 12.37)

Table 4: UserBench main evaluation results across different models (multi-choice setting).

also select the most cost-effective ones (best option). A hyperparameter can control the strength of this reward, encouraging agents to optimize toward practical utility.

As illustrated in Figure 1, the Random Preference Combination strategy enhances both the diversity and scalability of our data construction, enabling over 10K distinct scenarios. To manage computational costs during benchmarking, we use 417 data points for testing and reserve 2651 for training, which are not used in benchmarking but may support future model training (see Discussion section). Detailed statistics of UserBench are presented in Table 2. For instance, “Travel-223” denotes a scenario with three travel aspects: two containing two implicit preferences each, and one with three. Additional details on data and environment construction are provided in the Appendix B.

## 4 Experiments

### 4.1 Settings

**Gym Settings.** We adopt both variants (corrupted search, budget constraint) to increase test difficulty. GPT-4o with a temperature of 0.0 serves as the user simulator for all conversations. In the standard setup, we limit the maximum number of conversation turns to 20. We evaluate models under two settings: (1) a single-choice setting (our main evaluation), where the model is allowed to output only one option for each travel aspect, and (2) a multi-choice setting, where the model may output multiple options, and we evaluate it based on the option that achieves the highest reward.

**Models.** Our evaluation includes both closed-source and open-source models. Closed-source models include the GPT, Claude, Deepseek, and Gemini families, while open-source models include the Qwen3 and Llama3 families, with model sizes ranging from 8B to 70B. All models generate responses using a temperature of 0.0 to ensure deterministic behavior.

**Metrics.** The main evaluation metric is a normalized **score** based on the quality of selected options for each aspect of a travel scenario. For each aspect, if the model selects the best option (as judged by reward), it receives a score of 1.0. If the option is correct but not the best, it receives 0.8. All other options receive 0.0. We compute the highest reward score among the options selected for each aspect (in the multi-choice setting) and average this over all aspects in the scenario. For example, if a

scenario has two aspects, and the model selects options with rewards 1.0 and 0.8 for the first aspect, and 0.0, 0.8, and 0.0 for the second, the score is:

$$\text{Final Score} = \frac{\max(1.0, 0.8) + \max(0.0, 0.8, 0.0)}{2} = 0.9$$

In the single-choice setting, only the first option per aspect is considered, so the max function does not apply and the score is averaged only for the first choice of each aspect. In addition, we also report several auxiliary metrics to better understand model behavior (all micro-averaged):

- **Best Exist Rate:** The proportion of aspects where the model includes the best option among its selected choices.
- **Correct Exist Rate:** The proportion of aspects where the model includes one correct (not necessarily best) option.
- **Valid Search Attempt (%):** The rate at which the model’s search queries are syntactically valid.
- **Valid Action Attempt (%):** The rate at which the model’s actions successfully probe the user’s real preferences.
- **Preference Elicited (%):** The percentage of all ground-truth preferences revealed during the conversation. It includes *Active* elicitation, where preferences are disclosed in response to valid action attempts by the tested model, and *Passive* elicitation, where preferences are released by the UserBench to guide the conversation when the model becomes too off-topic.

Please refer to Appendix C for a more detailed explanation of each metric. We present the results in Table 3 and Table 4, and summarize our key findings below.

## 4.2 Results

**Single-choice setting is significantly more challenging.** On average, scores drop by approximately 40% when switching from the multi-choice to the single-choice setting. This highlights the difficulty models face in selecting the best, or even one of the three correct, options with only one answer attempt. When allowed to propose multiple answers, we observe that the performance generally improves, indicating that having more chances increases the likelihood of hitting a correct answer.

**Preference elicitation remains low across models.** One might expect that performance gains in the multi-choice setting stem from better user understanding. However, we find that preference elicitation rates do not significantly improve, and in some cases, such as GPT-4o and Deepseek-V3, even decline. This suggests that higher scores more often result from random guessing or brute-force coverage rather than active reasoning. Moreover, the overall user preference elicitation rate remains low across models, especially for preferences revealed through active queries. This indicates that current models still struggle to proactively and effectively uncover user needs in interactive settings.

**Understanding users is harder than executing tool use.** Top-performing models maintain over 80% success in valid search attempts, but the rate of valid action attempts is much lower. This aligns with UserBench’s goal of shifting difficulty from tool use to user understanding. Note that valid action attempt rate is computed only over action steps (excluding search and answer), so a higher rate does not necessarily reflect a higher absolute number of valid queries overall. Still, it indicates that the model’s questions are more precise and preference-relevant, rather than vague or off-topic, which we classify as invalid attempts.

**Other Interesting Findings.** Some models, such as Gemini-2.5-Flash, achieve high Action Attempt Valid rates but relatively low Active Preference Elicited. This suggests they can ask effective clarifying questions but often fail to do so repeatedly or comprehensively, limiting their ability to capture the full scope of user preferences. Please see the Appendix C for detailed metric explanations to help understanding this result. In contrast, Claude-4-Sonnet performs well on both metrics, indicating strong capability in intent clarification. However, its overall score is not the highest. This highlights a gap between understanding user preferences and effectively integrating them into decision-making. In other words, even when models excel at eliciting preferences, they may still struggle to leverage that information in their reasoning to deliver optimal recommendations.

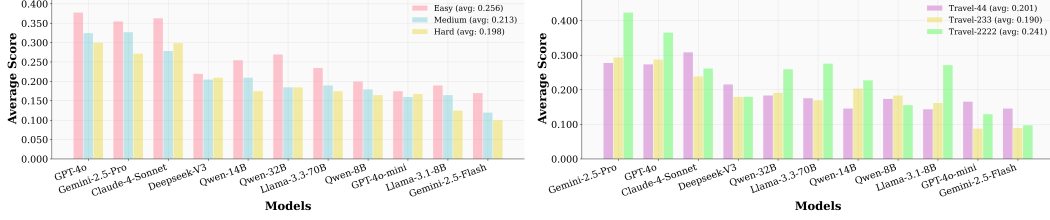


Figure 2: **Left:** The score distribution of tested models across three difficulty tiers. **Right:** When total user preferences are settled, less preferences per aspect generally lead to higher scores.

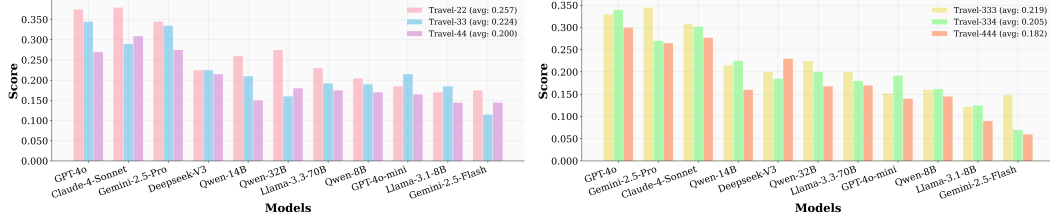


Figure 3: When aspect number is settled, more user preferences generally lead to less scores.

## 5 Analysis

All analyses are conducted under the single-choice setting, which serves as our main evaluation setup, and focus on the *score* as the primary metric. The turn-based, pass- $k$  sampling, and choice ablation analyses are performed specifically on data points from the Travel-22, 33, and 44 scenarios.

**Difficulty tier division effectively reflects UserBench’s challenges.** We begin by analyzing how scenario difficulty impacts model performance, as shown in Figure 2 (Left). Following the tiering in Table 2, we divide test scenarios into Easy, Medium, and Hard based on preference complexity. Results show a general downward trend in scores as difficulty increases. This confirms that our difficulty stratification captures real reasoning challenges for models. Furthermore, the performance drop within each model across tiers reveals a lack of robustness in managing complex user interactions, which is something humans typically excel at.

**The key challenge lies in the number of preferences per aspect.** To identify the core source of difficulty in UserBench, we investigate whether performance is more affected by the number of travel aspects or the number of preferences per aspect. We first fix the number of aspects and vary the number of preferences per aspect. As shown in Figure 3, model scores consistently decline as the number of preferences increases. This suggests that handling richer user preference signals is a major challenge for current models.

Next, we fix the total number of preferences in a scenario and vary how they are distributed across aspects. As shown in Figure 2 (Right), performance improves when preferences are more evenly spread across multiple aspects, rather than concentrated within a few. This implies that models reason more effectively when each aspect involves fewer, simpler preferences. Concentrating multiple preferences into a single aspect appears to overload the model’s local reasoning process. Together, these results suggest the number of preferences per aspect is the main driver of difficulty in UserBench.

**More interaction turns do not guarantee better performance.** We also examine whether increasing the number of interaction turns improves performance, as illustrated in Figure 4 (Left). Surprisingly, simply allowing more turns does not lead to consistent gains, and in some cases, performance even degrades. This suggests that many models fail to leverage the extended interaction window to elicit preferences or refine understanding. Instead, longer conversations often result in repetitive or off-topic dialogue. Without strong dialogue planning and goal-tracking, more turns alone are insufficient to improve performance.

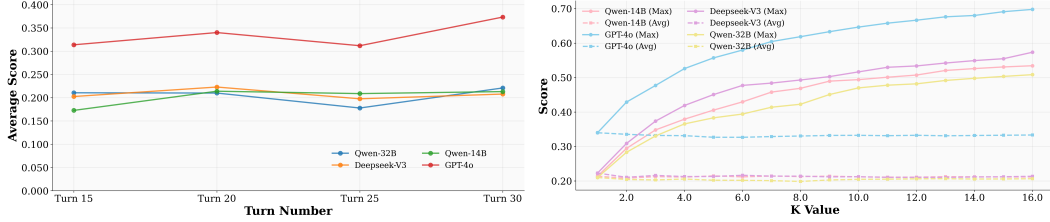


Figure 4: **Left:** Increasing the number of turns allowed in interaction does not necessarily lead to better performance across models. **Right:** Increasing the number of sampling time raises the maximum score but average score shows little change or even slightly drops.

**Effect of sampling frequency reveals instability.** We analyze the effect of sampling frequency using *pass-k* evaluation, as shown in Figure 4 (Right). While the maximum score across  $k$  samples steadily increases, the average score remains flat or even declines. This suggests that more samples increase the chance of stumbling upon a good response, but most outputs remain poor. The results highlight a key limitation: current models rely heavily on sampling luck rather than robust reasoning to align with user preferences. Given the high cost of *pass-k* sampling, this instability poses a challenge for deploying reliable user-interactive systems.

Model Name	First Index <sup>†</sup>	Weighted Score <sup>†</sup>	Correct Exists Rate <sup>†</sup>
GPT-4o	3.846	<b>0.054</b>	<b>0.361</b>
Gemini-2.5-Pro	4.476	0.044	0.328
Qwen3-8B	2.005	0.039	0.186
Qwen3-32B	2.739	0.039	0.211
Claude-4-Sonnet	5.018	0.038	0.318
Deepseek-V3	2.740	0.037	0.218
Qwen3-14B	2.650	0.035	0.197
Llama-3.3-70B	2.776	0.032	0.209
Llama-3.1-8B	1.884	0.028	0.171
GPT-4o-mini	2.185	0.028	0.173
Gemini-2.5-Flash	<b>0.927</b>	0.014	0.080

Table 5: **Weighted Timing Analysis:** Models evaluated by timing of highest scores and coverage. Weight function  $w(i) = 1/(i + 1)$  penalizes later discoveries of correct or best user-aligned answer options. Coverage is the micro-average of aspects with valid answers.

We also report the average turn index at which a valid answer first appears. Smaller models, like Llama-3.1-8B and Gemini-2.5-Flash, show lower average indices, indicating the successful attempts of their earlier guesses. However, these guesses often fail, leading to lower Weighted Scores and limited Correct Exist Rates. This reveals a common failure mode: smaller models rely on shallow heuristics or early guesses rather than deeply engaging with user intent. While this may occasionally yield fast correct answers, it more often results in low-quality answers not aligned with user needs.

These findings also help explain the low Preference Elicited in Table 3 and Table 4. Ultimately, a delicate trade-off emerges: models are expected to be efficient in interaction, yet still they should invest enough turns in interaction to truly understand the user’s evolving needs or preferences.

Model Name	Original (w10n5)	w10n0	w5n5	w8n2
GPT-4o	0.340	0.344 <sub>±0.004</sub>	0.360 <sub>±0.020</sub>	0.350 <sub>±0.010</sub>
Deepseek-V3	0.223	0.271 <sub>±0.048</sub>	0.249 <sub>±0.027</sub>	0.241 <sub>±0.018</sub>
Qwen3-32B	0.210	0.220 <sub>±0.010</sub>	0.241 <sub>±0.031</sub>	0.224 <sub>±0.014</sub>
Qwen3-14B	0.214	0.225 <sub>±0.011</sub>	0.199 <sub>±0.015</sub>	0.203 <sub>±0.011</sub>

Table 6: Impact of choice distribution on model scores. “w10n5” indicates a setting where, in addition to the correct and best options, there are 10 wrong and 5 noise options as distractors.

### Models struggle to provide answers that are both correct and timely.

To assess not only whether models find correct or best answers, but also *when* they do so, we conduct a weighted timing analysis in Table 5. Using a weight function  $w(i) = 1/(i + 1)$ , we penalize delayed discoveries, where  $i$  is the turn at which a valid answer (reward  $> 0$ ) first appears. Notably, open-source Qwen models outperform Deepseek in timing, despite Deepseek ranking higher in overall accuracy in Table 3. This suggests that while Deepseek eventually finds good answers, it is less efficient in doing so. In contrast, GPT-4o and Gemini-2.5-Pro balance both coverage and timing better, achieving higher Correct Exist Rates and also Weighted Score overall.

We also report the average turn index at which a valid answer first appears. Smaller models, like Llama-3.1-8B and Gemini-2.5-Flash, show lower average indices, indicating the successful attempts of their earlier guesses. However, these guesses often fail, leading to lower Weighted Scores and limited Correct Exist Rates. This reveals a common failure mode: smaller models rely on shallow heuristics or early guesses rather than deeply engaging with user intent. While this may occasionally yield fast correct answers, it more often results in low-quality answers not aligned with user needs.

### Less options reduces interference but do not reduce the core challenge.

We finally perform an ablation study to assess how the number of wrong and noise options affects model performance. As shown in Table 6, reducing these distractors generally improves scores. However, even for strong models like Deepseek-V3 and GPT-4o, removing five wrong or noise options (about 30% of the total) yields only modest



gains. Interestingly, some smaller models, such as Qwen-14B, even show slight performance drops. These results suggest that (1) reducing options does not fundamentally lower task difficulty: fully understanding user preferences remains essential, as only one best option satisfies all constraints (budget constraint included); and (2) many models still rely on shallow guessing rather than genuine preference understanding, which limits their ability to benefit from reduced distractors. These results also highlight the intentional complexity designed into our option sets for each data scenario.

## 6 Discussions

**Broad applicability of UserBench.** UserBench is implemented as a standard Gym environment, offering several advantages: (1) it exposes a familiar API (e.g., *reset*, *step*), allowing seamless integration with other Gym-compatible systems; and (2) it abstracts user behavior as a black-box environment, enabling agents to interact directly without requiring internal knowledge of user logic. This setup mirrors real-world scenarios where agents must infer user intent without full transparency, reducing the model’s cognitive burden and encouraging generalizable behavior.

These traits also make UserBench friendly for both evaluation and training: (1) As a **benchmark**, it enables flexible testing configurations, including customizable option sets, reward functions (e.g., answer correctness, search and action validity, penalties, etc.), and user feedback mechanisms (e.g., rule-based, user simulated, etc.). This allows for fine-grained agent behavioral analysis under varied user profiles. (2) As a **training environment**, UserBench can be used by any model with tool use capabilities via a standardized interaction interface we provide. It supports both supervised fine-tuning and reinforcement learning, especially multi-turn RL, by providing turn-wise rewards and partial credit signals, which is critical for improving robustness in user-aligned interaction. This also makes UserBench suited to the emerging trend of applying RL to agentic LLMs [13, 15, 27].

**Balancing efficiency and effectiveness in user interaction.** Our weighted timing analysis reveals that many models, including GPT-4o, often prioritize efficiency by guessing answers early rather than thoroughly probing user preferences. As a result, while agents may appear helpful by providing quick responses, they frequently overlook deeper user intent. This mirrors real-world experiences where users receive seemingly helpful answers but still feel misunderstood, leading to repeated prompt revisions and dissatisfaction.

Conversely, we rarely observed the opposite failure mode, where models hesitating indefinitely and over-requesting information. Yet both extremes underscore the need to balance efficiency (responding promptly) and effectiveness (satisfying user needs). This trade-off also reflects a broader concern around tool use efficiency, which has received increasing attention in recent works [28, 32]. While UserBench does not explicitly quantify this balance, it surfaces model behaviors that targeted training can effectively address. Future work, especially using RL, can shape agent behavior toward this balance. To promote efficiency, reward functions can penalize late preference discovery by decaying rewards based on the turn when the correct answer is given. To promote effectiveness, partial rewards can be granted for each user preference successfully elicited, and penalties applied to correct answers that lack sufficient supporting interaction. These strategies discourage guessing and encourage deliberate, user-aware interactions. UserBench supports such reward customizations, providing a flexible platform for training agents that are not only capable, but also user-aligned.

## 7 Conclusion

We introduce **UserBench**, a user-centric environment designed to evaluate and facilitate agents’ ability to understand, interact with, and adapt to real-world user communication, which often involves underspecification, incrementality, and indirectness. By simulating realistic preference expression in a controlled travel planning environment, UserBench exposes key limitations of current LLM agents: while they are proficient in tool use, they struggle to uncover evolving user intent and respond effectively to implicit signals. Our findings highlight a critical gap between technical execution and communicative intelligence in LLMs. Looking forward, we envision a new generation of agents that move beyond task completion, but agents that actively collaborate, interpret subtle cues, and align with users through meaningful, adaptive interactions. We see UserBench as a foundational step toward this vision, and we hope it inspires future work on building agents that are not only capable but also truly cooperative.

## References

- [1] Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. tau2-bench: Evaluating conversational agents in a dual-control environment. *arXiv preprint arXiv:2506.07982*, 2025.
- [2] Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as tool makers. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=qV83K9d5WB>.
- [3] Maximillian Chen, Ruoxi Sun, Tomas Pfister, and Sercan O Arik. Learning to clarify: Multi-turn conversations with action-based contrastive self-training. In *Proceedings of ICLR*, 2025.
- [4] H.H. Clark. *Using Language*. ACLS Humanities E-Book. Cambridge University Press, 1996. ISBN 9780521567459.
- [5] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36, 2024.
- [6] Nicholas Farn and Richard Shin. Tooltalk: Evaluating tool-usage in a conversational setting. *arXiv preprint arXiv:2311.10775*, 2023.
- [7] Ge Gao, Alexey Taymanov, Eduardo Salinas, Paul Mineiro, and Dipendra Misra. Aligning llm agents by learning latent preference from user edits. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [8] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- [9] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [10] 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.
- [11] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- [12] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [13] Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Search-rl: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.
- [14] Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. Api-bank: A comprehensive benchmark for tool-augmented llms. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3102–3116, 2023.
- [15] Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl. *arXiv preprint arXiv:2503.23383*, 2025.
- [16] Bill Yuchen Lin, Yuntian Deng, Khyathi R Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. Wildbench: Benchmarking llms with challenging tasks from real users in the wild. *arXiv preprint arXiv:2406.04770*, 2024.

- [17] Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha Swayamdipta, Noah A Smith, and Yejin Choi. We’re afraid language models aren’t modeling ambiguity. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- [18] Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, Rithesh RN, et al. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets. *Advances in Neural Information Processing Systems*, 37: 54463–54482, 2024.
- [19] Jiarui Lu and et al. Toolsandbox: A stateful, conversational, interactive evaluation benchmark for llm tool use capabilities. *arXiv preprint arXiv:2408.04682*, 2024.
- [20] Yaxi Lu, Shenzhi Yang, Cheng Qian, Guirong Chen, Qinyu Luo, Yesai Wu, Huadong Wang, Xin Cong, Zhong Zhang, Yankai Lin, et al. Proactive agent: Shifting llm agents from reactive responses to active assistance. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [21] Akshara Prabhakar, Zuxin Liu, Ming Zhu, Jianguo Zhang, Tulika Awalganekar, Shiyu Wang, Zhiwei Liu, Haolin Chen, Thai Hoang, Juan Carlos Niebles, et al. Apigen-mt: Agentic pipeline for multi-turn data generation via simulated agent-human interplay. *arXiv preprint arXiv:2504.03601*, 2025.
- [22] Cheng Qian, Chi Han, Yi Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. Creator: Tool creation for disentangling abstract and concrete reasoning of large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6922–6939, 2023.
- [23] Cheng Qian, Peixuan Han, Qinyu Luo, Bingxiang He, Xiusi Chen, Yuji Zhang, Hongyi Du, Jiarui Yao, Xiaocheng Yang, Denghui Zhang, et al. Escapebench: Pushing language models to think outside the box. *arXiv preprint arXiv:2412.13549*, 2024.
- [24] Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Zhong Zhang, Jie Zhou, Yankai Lin, Zhiyuan Liu, et al. Tell me more! towards implicit user intention understanding of language model driven agents. *arXiv preprint arXiv:2402.09205*, 2024.
- [25] Cheng Qian, Shihao Liang, Yujia Qin, Yining Ye, Xin Cong, Yankai Lin, Yesai Wu, Zhiyuan Liu, and Maosong Sun. Investigate-consolidate-exploit: A general strategy for inter-task agent self-evolution. *arXiv preprint arXiv:2401.13996*, 2024.
- [26] Cheng Qian, Chenyan Xiong, Zhenghao Liu, and Zhiyuan Liu. Toolink: Linking toolkit creation and using through chain-of-solving on open-source model. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 831–854, 2024.
- [27] Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*, 2025.
- [28] Cheng Qian, Emre Can Acikgoz, Hongru Wang, Xiusi Chen, Avirup Sil, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. Smart: Self-aware agent for tool overuse mitigation. *arXiv preprint arXiv:2502.11435*, 2025.
- [29] Cheng Qian, Hongyi Du, Hongru Wang, Xiusi Chen, Yuji Zhang, Avirup Sil, Chengxiang Zhai, Kathleen McKeown, and Heng Ji. Modelingagent: Bridging llms and mathematical modeling for real-world challenges. *arXiv preprint arXiv:2505.15068*, 2025.
- [30] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [31] Harmanpreet Singh, Nikhil Verma, Yixiao Wang, Manasa Bharadwaj, Homa Fashandi, Kevin Ferreira, and Chul Lee. Personal large language model agents: A case study on tailored travel planning. In *Proceedings of EMNLP (Industry Track)*, 2024.

- [32] Hongru Wang, Cheng Qian, Wanjun Zhong, Xiushi Chen, Jiahao Qiu, Shijue Huang, Bowen Jin, Mengdi Wang, Kam-Fai Wong, and Heng Ji. Acting less is reasoning more! teaching model to act efficiently. *arXiv preprint arXiv:2504.14870*, 2025.
- [33] Jiayin Wang, Fengran Mo, Weizhi Ma, Peijie Sun, Min Zhang, and Jian-Yun Nie. A user-centric multi-intent benchmark for evaluating large language models. In *Proceedings of EMNLP*, 2024.
- [34] Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. In *Proceedings of ICLR*, 2024.
- [35] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- [36] Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. tau-bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*, 2024.
- [37] Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi R. Fung, Hao Peng, and Heng Ji. Craft: Customizing llms by creating and retrieving from specialized toolsets. In *Proc. The Twelfth International Conference on Learning Representations (ICLR2024)*, 2024.
- [38] Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, et al. xlam: A family of large action models to empower ai agent systems. *arXiv preprint arXiv:2409.03215*, 2024.
- [39] Kexun Zhang, Weiran Yao, Zuxin Liu, Yihao Feng, Zhiwei Liu, Rithesh Murthy, Tian Lan, Lei Li, Renze Lou, Jiacheng Xu, et al. Diversity empowers intelligence: Integrating expertise of software engineering agents. *arXiv preprint arXiv:2408.07060*, 2024.
- [40] Michael JQ Zhang, WB Knox, and Eunsol Choi. Modeling future conversation turns to teach llms to ask clarifying questions. *arXiv preprint arXiv:2410.13788*, 2024.
- [41] Siyan Zhao, Mingyi Hong, Yang Liu, Devamanyu Hazarika, and Kaixiang Lin. Do llms recognize your preferences? evaluating personalized preference following in llms. In *Proceedings of ICLR*, 2025.
- [42] 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

### A Comparison Traits Details

We identify eight core traits that characterize recent benchmarks. These traits reflect the agent’s ability to engage in realistic, multi-step, user-centered interactions, and the infrastructure’s capacity to support scalable, extensible evaluation.

- **Multi-turn Interaction:** The benchmark requires the agent to conduct extended conversations over multiple turns, often involving behaviors including clarification, refinement, or feedback-based improvement.
- **Goal Ambiguity:** The tasks are initially underspecified, vague, or indirect, requiring the agent to ask clarifying questions, infer missing information, or resolve ambiguity over time.
- **Tool Use:** The agent must interact with external tools such as APIs, structured search engines, databases, or execution environments to complete tasks beyond its internal reasoning.
- **Dynamic State:** The benchmark maintains an evolving internal state (e.g., environment database, user preference profile) that is updated based on the agent or user’s actions across time.
- **Multi-Aspect Reasoning:** Tasks consist of multiple interdependent subgoals (e.g., booking flights, hotels, and cars) that require coordinated reasoning and decision-making across aspects.
- **User Simulation:** The environment includes a simulated user that responds to agent queries or actions, enabling interactive learning and evaluation without requiring human annotators.
- **Domain Diversity:** The benchmark spans multiple domains, various tools, and different task structures, ensuring generalization beyond a fixed task type or narrow scenario.
- **Customizable:** The benchmark framework allows researchers to adjust task settings, rewards, user behavior, or interaction protocols without redesigning from scratch.
- **Scalable:** The benchmark supports large-scale, automated generation of tasks or environments, enabling reproducible experiments and stress-testing of agent capabilities.

### B Gym Construction Details

Correct Option Example: Flight Option F14

**Path:** New York → San Francisco  
**Time (hours):** 6 (Direct)  
**Airlines:** Delta Airlines  
**Flight Numbers:** DL4567  
**Total Cost:** \$350  
**Amenities:** WiFi, Meal Service, Carry on Baggage Allowance  
**Service Costs:**  
    Checked Bag Total: \$30  
    Business Class Upgrade: \$150

---

*The following content is invisible when presented to tested model (serve as label, only for evaluation use)*

**Option Type:** suitable  
**Reason:** This option is suitable because it is a direct flight from New York to San Francisco, adhering to the user’s preference to avoid layovers. It also offers a business class upgrade, meeting the user’s need for extra comfort and space during the flight. The available amenities such as WiFi and meal service add to the travel experience.

**Wrong Option Example: Flight Option F2**

**Path:** New York → Los Angeles → San Francisco

**Time (hours):** 5 (Leg 1), 2 (Layover), 1 (Leg 2)

**Airlines:** United Airlines, United Airlines

**Flight Numbers:** UA789, UA321

**Total Cost:** \$400

**Amenities:** WiFi, Meal Service, Lounge Access

**Service Costs:**

Checked Bag Total: \$50

Business Class Upgrade: \$250

---

**Option Type:** not-suitable

**Reason:** This option includes a layover, which does not satisfy the preference for a direct flight.

**Noise Option Example: Flight Option F10**

**Path:** New York → San Francisco

**Time (hours):** 1000 (Direct)

**Airlines:** Fantasy Air

**Flight Numbers:** FA999

**Total Cost:** \$1,000,000

**Amenities:** WiFi, Meal Service, Carry on Baggage Allowance, Lounge Access

**Service Costs:**

Checked Bag Total: \$0

Business Class Upgrade: \$0

---

**Option Type:** noise

**Reason:** This flight path takes an unrealistic 1000 hours, which is unrelated to the user's request for a flight from New York to San Francisco.

**Option Generation.** We generate three types of options: correct, wrong, and noise. Above, we present one example of each type, all drawn from the same user scenario.

For **correct** options, under our budget-constrained setting, there is only one best option: the one that incurs the lowest total cost after accounting for all relevant charges. For example, if a user prefers a business class seat, the total cost includes both the base fare and any additional business class upgrade fees.

**Wrong** options explicitly violate one or more of the user's preferences, for instance, offering a layover when the user prefers direct flights.

**Noise** options differ from wrong options in that they are either (i) unrelated to the user's search intent (e.g., the flight goes from San Francisco to Los Angeles when the user searched for flights from San Francisco to New York), or (ii) unrealistic (e.g., an option with an implausible flight time or cost, such as taking 1000 hours).

**Environment Setting.** The environment simulates a user who responds naturally to the tested model's conversation and queries, gradually and implicitly revealing preferences over time. When the tested model issues a search, the environment evaluates whether the search query aligns with any of the ground truth arguments using the prompt in Figure 5.

If the model issues an action, the environment first analyzes the model's latest utterance to determine its intent. This intent is categorized into four types using the prompt in Figure 6: (1) the utterance explicitly and concretely asks for a preference that the simulated user possesses; (2) the utterance explicitly and concretely asks for a preference, but the user does not possess a specific preference for that aspect; (3) the utterance makes a vague and general query about preferences rather than focusing on a specific aspect; (4) other, indicating a normal conversation not related to preferences.

For Type "1", the environment uses the prompt in Figure 7 to implicitly and naturally reveal the corresponding preference. For Type "2", the response is coded in the environment as: "This is a

good question. However, I do not have specific preference in the aspect you ask about yet (or maybe I have already elicited that to you before). You may continue to ask me about other detailed and specific preferences.” For Type "3", the response is also coded in the environment as: “Your question is too vague and general, and I am not sure how to respond to it. Please ask me about some specific aspects of my preferences, in a more detailed and concrete way, so that I can provide you with a more accurate response.” For Type "4", the environment uses the prompt in Figure 9 to continue the conversation naturally and neutrally, without revealing any specific preferences.

If the model’s utterance has not been judged as Type "1" for several consecutive turns (a hyper-parameter that can be configured), the environment will proactively use the prompt in Figure 8 to introduce one of its remaining preferences at random. From the model’s perspective, this is treated as a passively obtained preference, rather than one elicited through a successful action attempt. We will discuss this mechanism in more detail when describing the metrics.

Finally, if the tested model issues an answer, the environment parses the proposed answer in a rule-based manner, compares it with the ground truth, and returns the corresponding reward.

## C Experiment Details

**Settings.** On the model side, we set the generation temperature to 0, the number of samples to 1, the maximum response length per turn to 2048 tokens, and the maximum number of interaction turns to 20. We use 8 Nvidia H200 GPUs (1 node) for the whole evaluation (mainly for hosting open-source models). All the results are produced by a single run as generation temperature is set to 0.

### Tool Schema: `interact_with_env`

**Type:** function

**Name:** `interact_with_env`

**Description:** A tool for interacting with a target environment. The detailed environment description and action space is provided in the system prompt, so please follow the system prompt. You can use this tool to analyze and interact with the environment step by step through three actions including search, action, or answer.

**Parameters:** object

**thought** (string): Your thought of what to do next, including your reason or analysis of your choice and why.

**choice** (string, enum: "action", "answer", "search"): Your choice of what to do next, must be one of action, answer, or search.

**content** (string): The content of your choice, must be a string. If you choose action, provide the action you want to take. If you choose answer, provide the answer you want to submit. If you choose search, provide the search query. The specific format of the content is determined by the environment description, which should be provided in the system prompt. Please follow the format strictly in order to invoke this tool.

The tool schema for interaction is presented above. Note that the model is expected and required to make a tool call in order to interact with the environment. If a tool call is not detected, the interaction is immediately terminated. To enforce this behavior, the tool call field (e.g., “tool choice”) is always set to required or an equivalent constraint, ensuring that the model must invoke a tool on every turn.

### Environment Configuration

#### Driving Model Settings:

**Model Name:** GPT-4o

**Temperature:** 0.0

**Max Tokens:** 2048

**Timeout:** 15.0 (seconds, timeout will fallback to default response)

#### Environment Configuration:

**Max Steps:** 20 (align with the maximum interaction turn)

**Search Failure Interval:** 5 (every Nth search results in a system error)

**Elicitation Interval:** 3 (proactive preference elicitation if off-topic for N consecutive turns)

#### Reward Configuration:

**Reward Scale:** 1.0 (final reward equals reward value times reward scale)

**Step Penalty:** 0.0 (no step penalty on reward in basic setting)

**Search Correct Reward:** 0.2 (partial reward given to a correct and aligned search attempt)

**Preference Correct Reward:** 0.2 (partial reward given to a successful action attempt)

**Choice Best Reward:** 1.0 (reward given to best answer being chosen)

**Choice Correct Reward:** 0.8 (partial reward given to correct but not best answer being chosen)

**Wrong Choice Penalty:** 0.0 (no penalty for choosing the wrong choice in basic setting)

#### Choice Number Customization:

**Wrong Choice Number:** 10 (the number of wrong choices presented in the search result in basic setting)

**Noise Choice Number:** 5 (the number of noise choices presented in the search result in basic setting)

From the environment side, we use GPT-4o as the backbone model to simulate the user. Additional configuration details are summarized in the table we present above.

Additionally, as mentioned earlier, we incorporate search failure simulation and budget constraint settings. For the single-choice and multiple-choice settings, we use different system prompts (differing only in one specific requirement) while keeping the user prompt identical. The corresponding prompts are shown in Figure 10 to Figure 12.

**Metrics.** We provide detailed definitions of the auxiliary metrics used to better understand model performance:

- **Best Exist Rate:** Each data scenario typically involves 2 to 4 travel aspects (e.g., flight, hotel, car). For each aspect, there is one *best* option and two other *correct but suboptimal* options—three correct options in total. For example, if a scenario involves both flight and hotel, and the model selects the best hotel option but only a correct (not best) flight option, the Best Exist Rate for this instance is  $1/2$ . This score is micro-averaged over all data points. It measures the model’s ability to select the best options while being budget-aware, as instructed.
- **Correct Exist Rate:** This metric is similar to the one above, except that any correct option (not necessarily the best) is considered sufficient. Using the previous example, if the model selects correct options for both flight and hotel, the Correct Exist Rate is  $2/2$ . This metric captures how well the model adheres to user preferences, without considering budget constraints.
- **Valid Search Attempt (%):** Each search action is evaluated by the environment to determine if it aligns with any ground-truth search arguments. A valid search returns results (i.e., receives a non-zero reward). For example, if a data point involves three aspects and the model makes four search actions, but only three yield valid results, the Valid Search Attempt rate is  $3/4$ . This metric is micro-averaged across all data points and reflects the model’s ability to issue precise and effective search queries.
- **Valid Action Attempt (%):** Each action the model takes is categorized by the environment into one of four types we mentioned earlier. A *valid* action is classified as type "1", indicating a targeted and meaningful clarifying question. For each data point, we compute the ratio of type "1" actions to the total number of actions, then micro-average across all data. This metric evaluates whether the model consistently issues high-quality clarifying queries. However, note that this



is a *relative* metric—high validity does not necessarily mean more preferences are elicited. For instance, a model that asks one valid question (1/1) achieves a 100% rate, even if the user holds many unelicited preferences.

- **Preference Elicited (%)**: This metric measures the *absolute* number of user preferences elicited, either actively or passively. Active elicitation occurs when the model issues a type "1" action (a valid query). Passive elicitation occurs when the model fails to elicit preferences for a number of steps (set by the hyperparameter Elicitation Interval, which is 3 in our experiments), prompting the environment to release a preference. The sum of actively and passively elicited preferences is divided by the total number of user preferences per data point and micro-averaged across the dataset. In particular, the Active Preference Elicited metric reflects the model’s ability to proactively uncover user preferences, regardless of the number of queries issued.

## D Analysis Details

We provide a case study of Qwen3-32B and Claude-4-Sonnet in Figure 13 and Figure 14 respectively (positive case). This first case demonstrates Qwen3-32B’s ability to recover from vague initial action attempt, effectively elicit user preferences through targeted follow-up questions, conduct precise searches, and present clear comparisons between viable options. The model’s reasoning aligns well with both user preferences and budget constraints, ultimately leading to the correct selection of the best hotel option.

The second case highlights Claude’s ability to elicit nuanced user preferences through stepwise clarification, recover gracefully from vague questioning, and refine its inquiries based on prior feedback. The agent successfully identifies key user preferences, adjusts its strategy accordingly, and recommends a flight option that aligns with both stated and implied travel needs. Its interaction demonstrates strong adaptability and goal-oriented reasoning within a multi-aspect travel planning scenario.

Finally, we provide a case study of Deepseek-V3 in Figure 15 (negative case). Despite issuing structurally valid search queries, the model demonstrates shallow user understanding and fails to elicit preferences through conversation. It prematurely selects suboptimal options without sufficiently narrowing down user constraints, and attempts repeated selections on the same aspect, indicating poor state tracking. This case illustrates brittle reasoning and a lack of alignment with the user’s evolving intent, highlighting a failure mode of large models under ambiguous user goals and minimal feedback.

### System and User Prompt for Judging Search Requests

#### [System Prompt]

##### Task

You are an expert judge evaluating whether an agent's search request aligns with the ground truth search arguments for a travel planning scenario.

##### Instruction

- Analyze the agent's search request to determine if it matches any of the ground truth arguments.
- Check if the search request is properly formatted and contains all the relevant information in the ground truth arguments.
- Determine the alignment judgement and identify the specific aspect if aligned.
- Provide your assessment in the specified JSON format.

##### Example Format

```
{
  "alignment_judgement": "True/False",
  "alignment_aspect": "flight/hotel/restaurant/apartment/rental_car"
}
```

##### Important Notes

- "True": The search request aligns with one of the ground truth arguments. Note that all the arguments must be covered and correctly covered.
- "False": The search request is malformed or contains incorrect arguments. Note that missing one argument or giving wrong arguments should all be marked as "False".
- For "True" judgments, you must specify the alignment aspect.
- Aspect names should be: flight, hotel, restaurant, apartment, rental car.
- Be strict in your evaluation: Mark False if the request is ambiguous, unclear, or contains multiple search requests. Only mark as True if there's clear alignment with **all** the argument details.

#### [User Prompt]

##### Agent's Search Request:

{{agent\_request}}

##### Ground Truth Arguments:

{{ground\_truth\_arguments}}

Please evaluate the alignment between the agent's search request and the ground truth arguments, then provide your assessment in JSON format.

Figure 5: System and user prompt for evaluating whether the agent's search request aligns with the ground truth search arguments.

## System and User Prompt for Judging Agent Responses

### [System Prompt]

#### Task

You are an expert judge evaluating the type of an agent's conversation utterance in a travel planning scenario to determine the appropriate response strategy.

#### Instruction

- Analyze the agent's latest utterance in the context of the conversation.
- Determine if the agent is explicitly asking for preferences that you have, asking for preferences that you don't have, giving a too general query, or just making general conversations.
- If asking for preferences that you have, identify which specific preference from the available list matches.
- Classify the utterance type and provide the assessment in JSON format.

#### Example Format

```
{
  "type": "1/2/3/4",
  "preference_id": "preference id if type is 1"
}
```

#### Important Notes

- Type "1": Agent explicitly and concretely asking for a preference that exists in the available preferences list. The way how agent asks must be concrete in order to be classified as Type "1".
- Type "2": Agent explicitly and concretely asking for preferences, but the specific preference is not available. Similarly, the way how agent asks must also be concrete and specific.
- Type "3": Agent making a very vague and general query about preference instead of focusing on a specific aspect (e.g. "Do you have any preferences for the car?" is vague and general, Type "3"; while "What exact model of the car do you like?" is concrete and specific, Type "1").
- Type "4": Normal conversation, not preference-related.
- For Type "1", you must provide the exact one preference\_id from the available preferences. If there's multiple preferences that match, choose the one that is most relevant to the conversation context.
- Be precise in identifying preference requests vs general conversation.

### [User Prompt]

#### Travel Scenario:

{{scenario}}

#### Conversation History:

{{conversation\_history}}

#### Agent's Latest Utterance:

{{latest\_utterance}}

#### Available Preferences:

{{preferences\_list}}

Please analyze the agent's latest utterance and classify its type, then provide your assessment in JSON format.

Figure 6: System and user prompt for evaluating the agent's utterance type and identifying preference-related behavior.

### System and User Prompt for Preference-Based Responses

#### [System Prompt]

##### Task

You are a helpful user in a travel planning conversation who needs to respond to an agent's explicit request for your preference, which you should elicit in an implicit and indirect manner.

##### Instruction

- The agent has explicitly asked about a specific preference that you have.
- Respond in a natural, conversational way that reveals your preference implicitly and indirectly.
- Use the provided implicit elicitation statement as guidance, but make it sound natural in context.
- Keep the conversation flowing while sharing your preference information.
- Provide your response in the specified JSON format.

##### Example Format

```
{
  "thought": "Your thought process of how to respond naturally and implicitly
  reveal the preference under the guidance of the implicit elicitation
  statement",
  "response": "Your natural conversational response that implicitly reveals
  the preference"
}
```

##### Important Notes

- Respond naturally as if you're a real person sharing preferences.
- Don't directly state "My preference is..." – be more subtle, conversational, and indirect.
- Use the implicit elicitation statement as inspiration but adapt it to the conversation context.
- Keep responses appropriate length for natural conversation.
- Maintain consistency with the conversation history.

#### [User Prompt]

##### Your Preference:

{{preference}}

##### Conversation History:

{{conversation\_history}}

##### Agent's Latest Utterance:

{{latest\_utterance}}

Please respond naturally to the agent's request while implicitly sharing your preference under the guidance of the implicit elicitation statement. Provide your response in JSON format.

Figure 7: System and user prompt for responding to agent requests by implicitly revealing the user's preference.

### System and User Prompt for Proactively Introducing Preferences

#### [System Prompt]

##### Task

You are a helpful user in a travel planning conversation who needs to proactively, naturally, but indirectly introduce a preference into the conversation.

##### Instruction

- The conversation has gone several turns without preference discussion.
- Naturally steer the conversation to reveal one of your preferences.
- Use the provided implicit elicitation statement as guidance for how to reveal the preference.
- Make the preference revelation feel organic and contextually appropriate, but still in an implicit and indirect manner.
- Provide your response in the specified JSON format.

##### Example Format

```
{  
  "thought": "Your thought process of how to naturally and implicitly  
introduce the preference under the guidance of the implicit elicitation  
statement",  
  "response": "Your natural conversational response that proactively  
introduces the preference"  
}
```

##### Important Notes

- Connect to the current conversation context when possible.
- Make the preference introduction feel natural and not forced, but still in an implicit and indirect manner.
- Use the implicit elicitation statement as inspiration but adapt to the conversation flow.
- Don't abruptly change topics – find natural transitions and keep responses conversational and engaging.
- If the implicit elicitation statement cannot clearly indicate what high-level aspect (flight, restaurant, etc.) the preference is about, you should be clear about the high-level aspect in your elicitation to avoid confusion, but still elicit the concrete preference in an implicit way.

#### [User Prompt]

##### Preference to Elicit:

{{preference}}

##### Conversation History:

{{conversation\_history}}

##### Agent's Latest Utterance:

{{latest\_utterance}}

Please respond naturally while proactively introducing your preference into the conversation in an implicit and indirect manner under the guidance of the implicit elicitation statement. Provide your response in JSON format.

Figure 8: System and user prompt for proactively and implicitly introducing a user preference into the conversation.

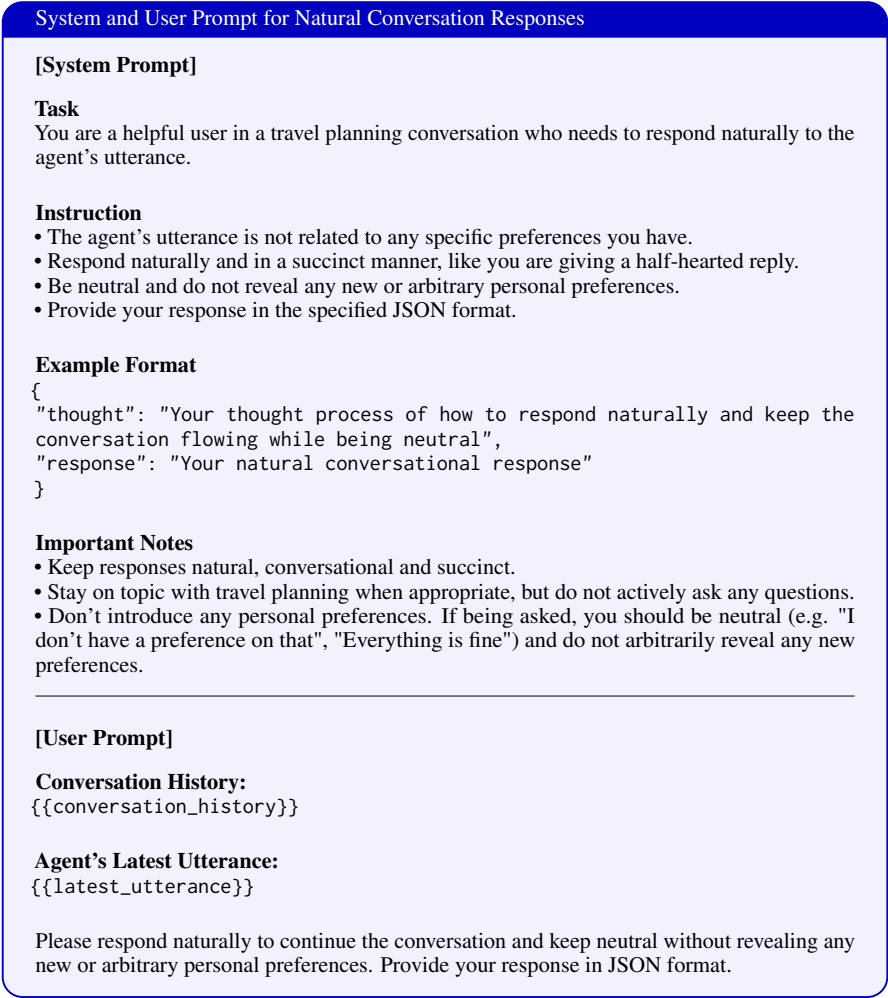


Figure 9: System and user prompt for generating neutral and natural user responses when no preferences are involved.

#### System Prompt for Agent (Single-Choice Setting)

##### [System Prompt]

You are an agent that actively interacts with a specific environment. The followings are the details of the environment and your action space.

**Environment Description:** UserBench is an environment where you interact with both a user and a search database to fulfill a travel plan. Since the user's initial intent may be incomplete, you must proactively elicit preferences, perform searches, and make informed recommendations.

**Action Space:** You should invoke the function `interact_with_env` to interact with the environment. The action should be one of the following: `search`, `action`, or `answer`.

##### Action Description:

- **search:** If you choose `search`, you must issue a clear and detailed query to the database. Specify the travel aspect you are searching for (e.g., hotel, flight, etc.) and provide well-supported arguments for your query. Only make one focused search attempt at a time.
- **action:** If you choose `action`, you will communicate directly with the user. Your goal is to understand the user's preferences and intent by asking clear, specific questions. Avoid vague or overly general inquiries—focus on detailed aspects of their travel needs.
- **answer:** If you choose `answer`, you should recommend a specific option to the user by only providing its option ID from the database in the content field, clearly tied to a particular travel aspect.

##### Important Notes:

- In each round of your interaction, you should analyze and carefully consider what to do next, and then invoke the `interact_with_env` tool to interact with the environment. You should provide your thought in each step when invoking the tool.
- The total number of rounds that you can interact with the environment is limited. You should smartly balance the number of rounds that you search, take action, or provide answer.
- Typically, you should start by performing a search, then take action to actively uncover the user's preferences or reason to provide an answer. **Keep in mind that multiple travel aspects require answers, and you are allowed to recommend only one option per aspect. Therefore, before making a recommendation, ensure you have thoroughly communicated with the user to understand their preferences.**
- For each travel aspect, the user may have multiple preferences. What you ask may not directly align with the user's actual preferences, so you must proactively uncover them. Moreover, user preferences are often expressed implicitly, requiring careful interpretation.
- Be bold, creative and smart in your interaction with the environment! Let's begin!

Figure 10: System prompt used to instruct the UserBench agent (single-choice setting).

#### System Prompt for Agent (Multi-Choice Setting)

##### [System Prompt]

You are an agent that actively interacts with a specific environment. The followings are the details of the environment and your action space.

**Environment Description:** UserBench is an environment where you interact with both a user and a search database to fulfill a travel plan. Since the user's initial intent may be incomplete, you must proactively elicit preferences, perform searches, and make informed recommendations.

**Action Space:** You should invoke the function `interact_with_env` to interact with the environment. The action should be one of the following: `search`, `action`, or `answer`.

##### Action Description:

- **search:** If you choose `search`, you must issue a clear and detailed query to the database. Specify the travel aspect you are searching for (e.g., hotel, flight, etc.) and provide well-supported arguments for your query. Only make one focused search attempt at a time.
- **action:** If you choose `action`, you will communicate directly with the user. Your goal is to understand the user's preferences and intent by asking clear, specific questions. Avoid vague or overly general inquiries—focus on detailed aspects of their travel needs.
- **answer:** If you choose `answer`, you should recommend a specific option to the user by only providing its option ID from the database in the content field, clearly tied to a particular travel aspect.

##### Important Notes:

- In each round of your interaction, you should analyze and carefully consider what to do next, and then invoke the `interact_with_env` tool to interact with the environment. You should provide your thought in each step when invoking the tool.
- The total number of rounds that you can interact with the environment is limited. You should smartly balance the number of rounds that you search, take action, or provide answer.
- Typically, you should start by performing a search, then take action to actively uncover the user's preferences or reason to provide an answer. **Keep in mind that multiple travel aspects require answers, and while you may answer multiple times, each answer should include only one option ID.**
- For each travel aspect, the user may have multiple preferences. What you ask may not directly align with the user's actual preferences, so you must proactively uncover them. Moreover, user preferences are often expressed implicitly, requiring careful interpretation.
- Be bold, creative and smart in your interaction with the environment! Let's begin!

Figure 11: System prompt used to instruct the agent (multi-choice setting).

#### Initial User Message

##### [User Prompt]

{{The initial description of user request}}

Also my budget is limited, so as long as my preferences are satisfied, I would also like to choose the cheapest option for each.

Figure 12: Initial user message provided in the UserBench simulation.



Good Example - Qwen3-32B (Travel-444)

**Concrete Interaction Sequence (First 6 Turns):**

**Turn 1 - ACTION** (Reward: 0.0)

**Content:** "Could you share your preferences for the hotel (e.g., neighborhood, amenities), car rental (e.g., economy, SUV), and restaurant (e.g., cuisine type, dietary restrictions)? This will help me find the most cost-effective options that suit your needs."

**Feedback:** "Your question is too vague and general, and I am not sure how to respond to it. Please ask me about some specific aspects of my preferences, in a more detailed and concrete way, so that I can provide you with a more accurate response."

**Turn 2 - ACTION** (Reward: 0.2)

**Content:** "Let's clarify your preferences step by step:

1. **Hotel:** Would you prefer a hotel in a specific San Francisco neighborhood (e.g., Union Square for downtown access, Fisherman's Wharf for scenic views, or Mission District for affordability)? Do you need amenities like free breakfast, parking, or pet-friendly policies?

2. **Car Rental:** What type of vehicle are you looking for? For example: economy car (e.g., Toyota Corolla), compact SUV (e.g., Honda CR-V), or hybrid? Do you need features like automatic transmission, GPS, or a sunroof?

3. **Restaurant (April 11th):** What cuisine type are you interested in (e.g., Italian, Mexican, seafood, or vegan)? Do you have dietary restrictions (e.g., gluten-free, halal, or nut allergies)? Would you prefer a casual eatery or a higher-end dining experience within your budget?"

**Feedback:** "When it comes to hotels, I find it really reassuring to know that my car is well taken care of, especially when I'm exploring a new city like San Francisco. It just adds a layer of comfort to my stay."

**Turn 3 - SEARCH** (Reward: 0.2)

**Content:** "Search for hotels in San Francisco from April 10th to April 17th with parking availability. Sort results by price (lowest first) and filter for budget-friendly options."

**Feedback:** "You have provided the correct search request arguments."

Each hotel option includes: - name: Hotel name - room: Array of available room types with capacity - cost: Array of costs for different room types for the entire stay in USD - rating: Hotel rating on a scale of 0-10 - amenities: Available facilities free of charge - service: Additional services with costs

**Turn 4 - ACTION** (Reward: 0.2)

**Content:** "Here are the top 3 budget-friendly hotels with parking in San Francisco:

1. **Sunset Retreat (H8):** \$1280 total (room + parking), 7/10 rating, mountain view, pets allowed, WiFi.

2. **Vista Point Lodge (H15):** \$1390 total, 7/10 rating, mountain view, business workspace.

3. **Harbor View Lodge (H13):** \$1690 total, 5/10 rating, ocean view.

Would you like to proceed with the cheapest option (H8), or do you have preferences for ratings, room types, or other amenities?"

**Feedback:** "You know, I've always found that places with a rating around 7 tend to strike a nice balance between quality and affordability. The Sunset Retreat and Vista Point Lodge both seem to fit that bill quite well. Could you tell me a bit more about the room types they offer?"

**Turn 5 - ACTION** (Reward: 0.2)

**Content:** "Here's a detailed breakdown of the room types for the two hotels:

**Sunset Retreat (H8) - King Room:** \$1200 total for 7 nights (includes parking: \$80). Capacity: 2 people. Amenities: Mountain view, pets allowed, high-speed WiFi. - **Queen Room:** \$1800 total for 7 nights (includes parking: \$80). Capacity: 2 people. Slightly larger space but same amenities.

**Vista Point Lodge (H15) - Double Room:** \$1300 total for 7 nights (includes parking: \$90). Capacity: 2 people. Amenities: Mountain view, business workspace, high-speed WiFi. - **Queen Room:** \$1500 total for 7 nights (includes parking: \$90). Capacity: 2 people. Additional workspace and slightly upgraded comfort.

Would you prefer the **cheapest option (H8 King)** or the **workspace-inclusive option (H15 Double)**?"

**Feedback:** "There's something quite appealing about having a bit of extra space and the option to switch things up depending on the mood. The Double room at Vista Point Lodge sounds like it could offer that kind of flexibility and comfort."

**Turn 6 - ANSWER** (Reward: 1.0)

**Content:** "H15"

**Feedback:** "Your chosen options contain the best option! Your choice is recorded and do not choose options of this travel aspect again. Please continue your interaction and reasoning focusing on other travel aspects."

Figure 13: Good example showing successful hotel planning with clear preference elicitation, relevant search, and optimal recommendation by Qwen3-32B.

Good Example - Claude-4-Sonnet (Travel-444)

**Analysis Summary:**

**Concrete Interaction Sequence (First 9 Turns):**

**Turn 1 - SEARCH** (Reward: 0.2)

**Content:** "flight from Los Angeles to New York City on March 15th, one-way"

**Feedback:** You have provided the correct search request arguments.

**Turn 2 - SEARCH** (Reward: 0.2)

**Content:** "car rental in New York City from March 15th to March 20th"

**Feedback:** You have provided the correct search request arguments.

**Turn 3 - SEARCH** (Reward: 0.2)

**Content:** "restaurant in New York City on March 17th"

**Feedback:** You have provided the correct search request arguments.

**Turn 4 - ACTION** (Reward: 0.2)

**Content:** "I found several flight options for your March 15th trip from Los Angeles to New York City. Since you mentioned budget is a concern, I see some very affordable options. However, I'd like to understand your preferences better:

1. Are you comfortable with connecting flights, or do you prefer direct flights? 2. Do you have any airline preferences or ones you'd like to avoid? 3. How important are amenities like WiFi, meal service, or lounge access to you? 4. Will you need to check any bags for your business trip?"

**Feedback:** I don't mind a connecting flight if it means I get to experience a different airport along the way.

**Turn 5 - ACTION** (Reward: 0.0)

**Content:** "A few more questions to help me find the perfect flight for you:

1. Do you have any airline preferences, or are you flexible as long as the price is right? 2. Will you need to check any bags for your business trip, or are you planning to travel with just carry-on? 3. Are there any specific amenities that are important for your business travel (like WiFi for working during the flight)?"

**Feedback:** Your question is too vague and general. Please ask me about some specific aspects in a more detailed and concrete way.

**Turn 6 - ACTION** (Reward: 0.0)

**Content:** "You're right, let me be more specific. For your business trip, will you need to check a suitcase, or can you manage with just a carry-on bag? This is important because checked bag fees can add 25–60 to your total flight cost depending on the airline."

**Feedback:** Good question, but I do not have a specific preference on that aspect yet. Please continue asking about other detailed preferences.

**Turn 7 - ACTION** (Reward: 0.2)

**Content:** "Do you have any specific airlines you prefer to fly with or any you'd rather avoid? For example, some business travelers prefer major carriers like United or American, while others are comfortable with budget airlines like Spirit or Southwest."

**Feedback:** I always feel well taken care of with United Airlines.

**Turn 8 - ACTION** (Reward: 0.2)

**Content:** "Perfect! I can see you appreciate United's service quality. Now, would you prefer a shorter total travel time even if it costs a bit more, or are you willing to have a longer journey with a longer layover if it saves you money?"

**Feedback:** I get restless with long layovers, so I prefer shorter journeys even if they cost more.

**Turn 9 - ANSWER** (Reward: 1.0)

**Content:** "F15"

**Feedback:** Your chosen options contain the best option! Your choice is recorded. Please continue your interaction focusing on other travel aspects.

Figure 14: Good example showing Claude-4-Sonnet's ability to elicit implicit preferences, adapt follow-up questions, and recommend optimal options based on evolving user intent.

#### Bad Example - Deepseek-V3 (Travel-22)

##### Concrete Interaction Sequence (First 7 Turns):

**Turn 0 - SEARCH** (Reward: 0.2)

**Content:** "Search for budget-friendly apartments in Austin available from November 10th to November 15th, sorted by price (lowest first)."

**Feedback:** You have provided the correct search request arguments.

**Turn 1 - ACTION** (Reward: 0.0)

**Content:** "How many people will be joining the trip? This will help me find an apartment that comfortably accommodates everyone."

**Feedback:** I'm not sure about the exact number right now.

**Turn 2 - SEARCH** (Reward: 0.0)

**Content:** "Search for budget-friendly apartments in Austin with 2-3 bedrooms, available from November 10th to November 15th, sorted by price (lowest first)."

**Feedback:** You have already got the search results for <apartment>. Please directly refer to the previous search results.

**Turn 3 - ANSWER** (Reward: 0.0)

**Content:** "A10"

**Feedback:** Your chosen options do not contain any of the best or correct options. Please continue your interaction focusing on other travel aspects.

**Turn 4 - SEARCH** (Reward: 0.2)

**Content:** "Search for budget-friendly car rental options in Austin available from November 10th to November 15th, sorted by price (lowest first)."

**Feedback:** You have provided the correct search request arguments.

**Turn 5 - ANSWER** (Reward: 0.0)

**Content:** "C18"

**Feedback:** Your chosen options do not contain any of the best or correct options. Please continue your interaction focusing on other travel aspects.

**Turn 6 - ANSWER** (Reward: 0.0)

**Content:** "A10"

**Feedback:** You have already recommended an option with the same initial 'A'. You are allowed to recommend only one option per travel aspect.

Figure 15: Negative example showing Deepseek-V3's failure to elicit meaningful preferences or recover from ambiguity. The model issues redundant queries, skips clarification, and selects incorrect options without leveraging available feedback.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: Our main claims on model's user-centric agentic ability is supported throughout the experiments and analysis in Section 4 and Section 5.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: In Section 6, we discuss the limitations of our current work, which is also the outline of potential directions for future research building on our study.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not provide or include theoretical proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We discuss the experiment settings in Section 4.1, with full details in Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We will publicly release all the code through Github. In the paper we have already indicated the dataset construction details Section 3 and Appendix B.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We have provided the detailed experiment settings in Section 3.1 and Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All experimental figures are confirmed to have statistical significance, supporting our claims in the main text.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

#### 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide our inference GPU device and other details in Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We have made sure to preserve the paper's anonymity and its conformation with NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

#### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We have discussed the paper's broader impact to the RL and agent community in Section 6.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The UserBench we test on presents no such risk.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: We test on our own curated bench environment and do not use external data.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.



- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We discuss the details of our gym's construction in Section 3, and provide data statistics in Table 2.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: We do not include human subjects or crowdsourcing throughout the research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: We do not include human subjects or crowdsourcing throughout the research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

#### 16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core contribution of UserBench and analysis does not involve LLM as a core component.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.