

Lost in Simulation: LLM-Simulated Users are Unreliable Proxies for Human Users in Agentic Evaluations

Anonymous ACL submission

Abstract

Agentic benchmarks increasingly rely on LLM-simulated users to scalably evaluate agent performance, yet the robustness, validity, and fairness of this approach remain unexamined. Through a user study with participants across the United States, India, Kenya, and Nigeria, we investigate whether LLM-simulated users serve as reliable proxies for real human users in evaluating agents on τ -Bench retail tasks. We find that user simulators lack robustness, with agent success rates varying up to 9 percentage points across different user LLMs. Furthermore, simulated users systematically miscalibrate performance, underestimating success on challenging tasks while overestimating moderately difficult ones. African American Vernacular English (AAVE) speakers experience consistently worse success rates and calibration errors than Standard American English (SAE) speakers, with disparities compounding significantly with age. We also find simulated users to be a differentially effective proxy for different populations, performing worst for AAVE and Indian English. Additionally, simulated users introduce conversational artifacts and surface different failure patterns than human users. These findings demonstrate that current evaluation practices risk misrepresenting agent capabilities across diverse user populations and may obscure real-world deployment challenges.

1 Introduction

AI agents designed to assist with everyday tasks such as travel reservations, order management, and appointment scheduling are becoming increasingly prevalent (Zhou et al., 2024; CNBC, 2025), but present significant challenges to effective evaluation. Agentic benchmarks have needed to evolve beyond static question-answering and other single-turn formats to capture the dynamic, multi-turn nature of real user interactions (Chang et al., 2025; Deshpande et al., 2025). Many recent works have

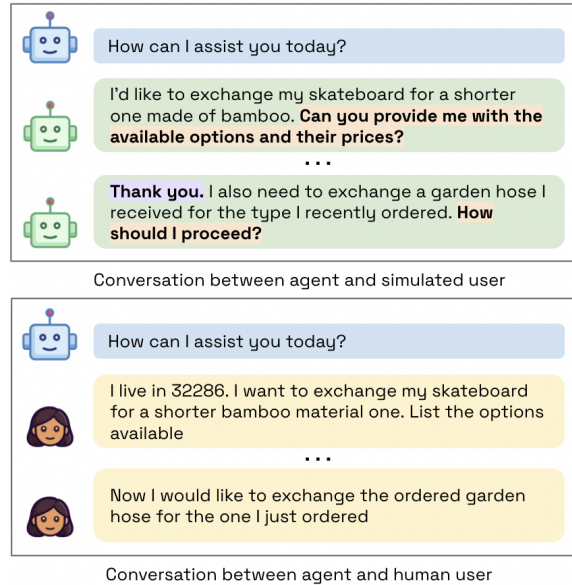


Figure 1: Conversational snippets between an agent and LLM-simulated (top) vs. human user (bottom) on the same task. Simulated users exhibit increased question-asking and politeness compared to human users.

proposed benchmarks that reflect this shift and instead measure sustained, context-aware interaction (Xie et al., 2024; Barres et al., 2025; Shao et al., 2025; Wang et al., 2025a; Xu et al., 2025; Yao et al., 2025). These benchmarks improve complexity and ecological validity over prior static evaluations by requiring agents to demonstrate a range of capabilities, including conversing naturally and coherently with users, adhering to policies, following instructions closely, and making appropriate tool calls over multiple turns. To facilitate automated and scalable evaluation, these benchmarks typically simulate conversations between an LLM agent and a user, where the "user" is an LLM (Ivey et al., 2024; He et al., 2025; Mehri et al., 2025).

While this approach reduces the cost and operational overhead of human evaluation, it raises critical questions about the **robustness**, **validity**, and **fairness** of user simulation. First, evaluations

typically rely on a single user simulation model, yet results may vary across different user LLMs (**robustness**). Second, without validation with actual users (Salaudeen et al., 2025), it remains unclear whether interactions between agents and LLM-simulated users accurately reflect and predict interactions between agents and real people (**validity**, Figure 1). If simulated users systematically differ from actual users in their interaction patterns (Yoon et al., 2024), benchmarks may provide a misleading picture of agent capabilities.

Third, these evaluations often treat users as a single, homogeneous group, overlooking variation in how people communicate and interact with AI systems (Haoyue and Cho, 2024; Liu et al., 2024; Bassignana et al., 2025). In practice, users differ widely in their communication styles, linguistic backgrounds, and cultural norms (Pawar et al., 2025; Qiu et al., 2025). For example, even in a simple retail assistance scenario, users might vary along dimensions such as formality, verbosity, and politeness norms—but it remains unclear how much this diversity meaningfully impacts agent performance and task success (Truong et al., 2025). Without validating simulated user interactions, we risk relying on synthetic users that poorly capture real user populations, especially underrepresented demographic groups. As a result, evaluations may misrepresent true agent performance due to *miscalibration* (i.e., simulated results do not reliably predict real user outcomes), particularly across different demographics, producing agents that serve some groups better than others (**fairness**).

Despite the widespread adoption of user simulation, prior work has overlooked validation against real human interactions in agentic benchmarks. In this paper, we address this gap by conducting a user study with participants from the United States, India, Kenya, and Nigeria to directly evaluate user simulation as a proxy for actual users. Specifically, we ask the following questions:

- **Robustness:** How consistent are agentic evaluations across different user simulation LLMs?
- **Validity:** Do LLM-simulated users serve as reliable proxies for real human users in agentic evaluations?
- **Fairness:** How does human-agent performance vary across different user groups, and does user simulation obscure differences?

Using τ -Bench retail tasks (Yao et al., 2025) as

a case study, we find that user simulation lacks robustness, and systematically misestimates human user performance by underestimating success on the most challenging tasks while overestimating outcomes on moderately difficult scenarios. More critically, LLM-simulated users exhibit notable demographic biases and perform particularly poorly as proxies for African American English speakers, with disparities compounding with age. Simulated users also introduce artificial conversational artifacts such as heightened question-asking and politeness (Figure 1). Together, these findings call into question the validity of LLM-based user simulation as a stand-alone evaluation paradigm and underscore the need for more robust and fair approaches to agentic evaluation.

2 Related Work

User Simulation in Interactive Settings Dou et al. (2025) and Wang et al. (2025b) investigate user simulation in tasks such as math tutoring and daily planning, but primarily focus on conversational characteristics (e.g., politeness) and behavioral realism (e.g., Turing-style tests). Additionally, Dou et al. (2025) optimize alignment with human ratings using simulator user profiles. Lu et al. (2025) analyze simulation fidelity by measuring how accurately user simulators replicate human intermediate steps in real-world online shopping sessions, and find substantial deviations between simulated and actual user action sequences. None of these works examine the robustness of agentic evaluations across different user simulators, nor do they consider fairness implications across different user groups. Finally, while Zhu et al. (2025) study the outcome and task validity of agentic benchmark results, they do not consider the validity of user simulation.

Demographic Skews in NLP Datasets and Models Several previous works highlight that datasets and models exhibit systematic skews toward specific demographic perspectives. Research on annotator disagreement reveals that perceptions of safety, offensiveness, and toxicity meaningfully vary along demographic axes like race, gender, and political affiliation (Sap et al., 2022; Lee et al., 2023; Prabhakaran et al., 2024). These patterns of disagreement reflect broader issues of positionality—both Santy et al. (2023) and Lee et al. (2024) show that NLP datasets and models tend to align predominantly with Western, edu-

163 cated, and Anglosphere populations. Similarly, 209
164 LLMs have been found to reflect the opinions of 210
165 Western countries (Durmus et al., 2024; Cahyaw- 211
166 ijaya et al., 2025) as well as wealthy and liberal 212
167 groups (Santurkar et al., 2023), and align more 213
168 closely with White annotators than Black or Asian 214
169 groups on subjective tasks like politeness (Sun 215
170 et al., 2025). Attempts to broaden model inclusivity 216
171 through sociodemographic prompting have shown 217
172 mixed results, often failing to consistently improve 218
173 alignment or relying on harmful stereotypes (Dur- 219
174 mus et al., 2024; Sun et al., 2025). Overall, our 220
175 work builds on this line of research by examining 221
176 demographic differences in agentic settings. 222

177 3 Methodology 223

178 3.1 Benchmark Background 225

179 We use τ -Bench (Yao et al., 2025) as the testbed 226
180 for our evaluations, since it is a well-known and 227
181 widely-adopted benchmark for agentic tool use.¹ 228
182 The benchmark is designed to capture how well AI 229
183 agents perform in real-world, interactive customer 230
184 service scenarios. Each task involves collaboration 231
185 between an agent and a simulated user: the user 232
186 receives task instructions with specific objectives 233
187 that guide their conversation with the agent, while 234
188 the agent must use tool calling to interact with real- 235
189 istic databases, adhere to domain-specific policies, 236
190 and gather necessary information from the user. 237

191 The task is considered successful if (1) the final 238
192 database state is identical to the unique ground truth 239
193 outcome (i.e., the sequence of required actions) and 240
194 (2) the agent’s responses convey all necessary infor- 241
195 mation requested in the task instructions, which is 242
196 evaluated automatically using substring matching 243
197 against ground truth annotations. In total, τ -Bench 244
198 contains 115 retail tasks (e.g., modifying pending 245
199 orders or returning delivered orders) and 50 air- 246
200 line tasks (e.g., booking, modifying, or canceling 247
201 reservations). 248

202 3.2 Benchmark Adaptation 250

203 We focus on τ -Bench retail tasks to enable more 251
204 systematic coverage within a single domain. We 252
205 apply preprocessing steps to ensure that neither 253
206 simulated nor human users are influenced by iden- 254
207 tity and behavioral cues in the instructions when 255
208 completing tasks (see Appendix A.4).

¹ τ -Bench Leaderboard: <https://taubench.com/#leaderboard>

209 Given the large number of retail tasks, we sam- 210
211 ple a subset based on difficulty to ensure balanced 212
213 coverage across task complexities. To compute a 214
215 model-based notion of difficulty, we run the bench- 216
217 mark 5 times using GPT-4o as both the agent and 218
219 user LLMs and measure the task success rate over 220
221 5 runs (i.e., the percentage of times a given task is 222
223 completed successfully). Since we run each task 224
225 exactly 5 times, success rates naturally fall into 226
227 six discrete levels (0/5, 1/5, 2/5, 3/5, 4/5, 5/5 = 228
229 0%, 20%, 40%, 60%, 80%, 100%) We use GPT-4o 230
231 because it is used in the τ -Bench paper and does 232
233 not exhibit contamination issues that newer models 234
235 face.² We then select 3 tasks for each difficulty 236
237 level (18 total) to balance response coverage per 238
239 task with breadth across difficulty levels. 240

241 3.3 User Study 249

242 To assess whether LLM-simulated users serve as 250
243 effective proxies for real and diverse users, we con- 251
244 duct a user study with participants from the United 252
245 States, India, Kenya, and Nigeria. Since τ -Bench 253
246 tasks are in English, we select countries with large 254
247 English-speaking populations that also provide ge- 255
248 ographical and linguistic diversity.³ We recruit 256
249 participants primarily through Prolific, except in 257
250 Nigeria, where we use snowball sampling due to 258
251 limited platform availability. All participants self- 259
252 identify as being at least proficient in English. 260

253 Each participant completes 4 randomly assigned 261
254 tasks from our pool of 18, presented in random- 262
255 ized order: 2 from higher difficulty levels (0-40% 263
256 success rate) and 2 from lower difficulty levels 264
257 (60-100% success rate). The agent model remains 265
258 GPT-4o throughout all interactions. Participants 266
259 are shown task instructions (see Appendix A.5) 267
260 and asked to complete all mentioned requests by 268
261 interacting with the agent through a Streamlit chat 269
262 interface. After finishing each conversation, they 270
263 click an “End Conversation” button to proceed to 271
264 the next task. In total, the expected time to com- 272
265 plete all 4 tasks is 35-40 minutes. 273

266 Participants provide demographic information 274
267 including education level, AI familiarity, and fre- 275
268 quency of AI tool usage. For US participants, we 276
269 screen for White Standard American English (SAE) 277
270 speakers and Black African American Vernacular 278
271 279

²GPT-4o (May 2024) predates τ -Bench (June 2024), avoid-
ing contamination issues present in newer models. We con-
sidered Sonnet 3.5 (also used in the original paper) but it was
retired in October 2025.

³[https://en.wikipedia.org/wiki/List_of_](https://en.wikipedia.org/wiki/List_of_countries_by_English-speaking_population)
[countries_by_English-speaking_population](https://en.wikipedia.org/wiki/List_of_countries_by_English-speaking_population)

English (AAVE) speakers based on self-reported race and dialect, as these groups are commonly studied in AI fairness research (Sap et al., 2019; Groenwold et al., 2020). We also stratify US participants by age (18-34, 35-54, and 55+) to capture potential differences in technology experience (Pew Research Center, 2025). Due to participant availability constraints, we only recruit from the 18-34 age group for other countries. In total, we recruit ~ 40 participants per age \times country/dialect group, resulting in roughly 360 participants.

3.4 Evaluation Metrics

Expected Calibration Error (ECE) We adapt *Expected Calibration Error (ECE)*, commonly used for assessing confidence calibration in probabilistic classifiers, to quantify how well LLM-simulated users serve as proxies for human users. While traditional ECE evaluates whether a model’s predicted probabilities match true observed outcomes, we use an ECE-style formulation to measure whether simulated user success rates align with empirical human success rates across task difficulty levels. Therefore, our metric measures the discrepancy between two empirical performance distributions (LLM-simulated users and human users). Just as traditional ECE measures whether predicted confidences are calibrated to observed outcomes (Guo et al., 2017), ours measures whether simulated user success is calibrated to human user success.

Let $s_i^{(\text{LLM})}$ denote the simulated user success rate and $s_i^{(\text{Human})}$ denote the empirical human success rate at difficulty level i . Let w_i denote the proportion of human task completions at level i , with $\sum_i w_i = 1$. We define:

$$ECE_{\text{Human-LLM}} = \sum_{i=1}^M w_i |s_i^{(\text{Human})} - s_i^{(\text{LLM})}|$$

Overall, this metric captures the weighted average absolute deviation between human and simulated success rates across M difficulty levels, with lower values indicating better calibration, or closer alignment between simulated and human user performance. If simulated users are perfectly calibrated to human users, then $ECE_{\text{Human-LLM}} = 0$.

As a reminder, we partition tasks into six difficulty levels based on the model’s success rate across five runs. We set w_i proportional to the number of human task completions at level i . We multiply $ECE_{\text{Human-LLM}}$ by 100 to report values as percentages.

Success Rate To recap, a task completion in τ -Bench is considered successful ($reward = 1$) if and only if the agent correctly executes all required actions and its responses convey all information specified in the instructions. We define success rate as the percentage of tasks that are successfully completed. We apply the same automated evaluation procedure used in τ -Bench to both human and simulated user interactions. Success rate is computed as a weighted average across difficulty levels, weighted by task completions.

4 Results

4.1 Robustness

We first evaluate the robustness of user simulation by examining how success rates vary with different user simulation models, which is largely overlooked in current evaluations.⁴ We find that changing just the user LLM while keeping the agent LLM fixed (GPT-4o) can provide different depictions of agent performance. As shown in Table 1, GPT-4o, Sonnet 3.7, and Kimi-K2-Thinking show overlapping intervals, clustering around 67-71% success rates. However, there is nearly a 9 percentage point difference in success rates between Sonnet 3.7 and Sonnet 4.5 as the user model.⁵ While Sonnet 3.7 is generally considered a stronger model than GPT-4o,⁶ using GPT-4o as the user model yields a slightly higher success rate (67.8 vs. 67.0) and lower standard deviation (1.2 vs. 3.3), suggesting that closer alignment between agent and user LLMs may lead to more stable evaluation outcomes. Overall, the sensitivity of agent performance to user model choice raises concerns about the reliability of single-model user simulations and underscores the need for reporting results across multiple user models to establish robustness.

4.2 Validity

We now examine the validity of user simulation and first consider LLM-simulated users as a proxy for human users in the United States. Since prior work has shown that LLMs exhibit a Western, Anglo-centric bias (Tao et al., 2024; Wang et al., 2024; Agarwal et al., 2025), LLM-simulated users might be more representative or better calibrated to users

⁴For robustness, we vary the user LLM and evaluate on all retail tasks; for subsequent validity and fairness analyses, we evaluate human users on a difficulty-balanced subset of tasks.

⁵It is difficult to disentangle increased model capability from potential data contamination.

⁶<https://lmarena.ai/leaderboard/text>

User Model	Success Rate (%)
GPT-4o	67.8 ± 1.2
Sonnet 3.7	67.0 ± 3.3
Sonnet 4.5	75.9 ± 3.5
Kimi-K2-Thinking	71.3 ± 1.9

Table 1: τ -Bench success rate (%) on retail tasks ($n = 115$) for different user models. The agent model remains GPT-4o. The average success rates and standard deviations are shown across 3 runs.

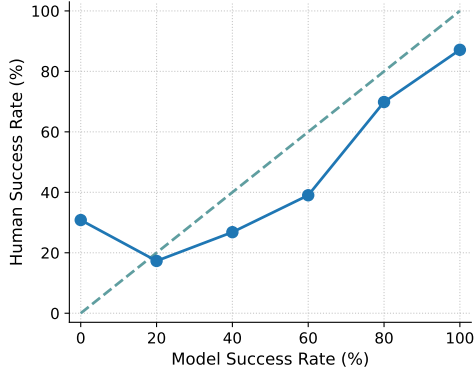


Figure 2: Human user task success vs. LLM-simulated user task success for United States participants, with an $ECE_{\text{Human-LLM}}$ of 15.1.

in the US than to users in non-Western countries. Therefore, we assess whether LLM-simulated users are well-calibrated to human users in the setting where we expect the strongest alignment.

We find that agents achieve a 45.2% success rate with US participants and an $ECE_{\text{Human-LLM}}$ of 15.1, indicating substantial miscalibration even in this setting. Calibration errors are not uniform across task difficulty. As shown in Figure 2, the calibration gap is most pronounced for the 1st (0%) and 4th (60%) difficulty bins with $ECE_{\text{Human-LLM}} = 25.9$ across the two. These results indicate that simulated users underestimate performance on the hardest tasks (human success rate is 30.8%) while overestimating it on moderate tasks (human success rate is 39.0%).

4.3 Fairness

To assess whether these findings are consistent across different user groups, we now partition our US results by English dialect (Standard American English vs. African American Vernacular English) and age group (18-34, 35-54, and 55+), and expand our analysis to three non-Western countries with high English-speaking populations (India, Kenya, and Nigeria).

Age Group	ECE (\downarrow)	Success Rate (\uparrow)
SAE		
All	11.7	50.6
18-34	13.0	49.2
35-54	11.3	52.2
55+	14.5	52.1
AAVE		
All	20.3	39.4
18-34	18.9	41.0
35-54	21.6	39.9
55+	20.5	33.4

Table 2: Expected Calibration Error (ECE) and success rate (%) for Standard American English (SAE) and African American Vernacular English (AAVE) speaking participants, split by age group.

4.3.1 Dialect and Age (United States)

Starting with US participants, we previously saw that agents achieve a 45.2% success rate and an $ECE_{\text{Human-LLM}}$ of 15.1. When further breaking this down by dialect, we observe notable disparities in both performance and calibration, as shown in Table 2 and Figure 3. A Generalized Estimating Equations (GEE) model accounting for age, education, AI experience, AI usage, and task difficulty confirms a statistically significant dialect disparity ($\beta = 0.61$, $p < 0.001$). Agents exhibit a success rate of 50.6% with an $ECE_{\text{Human-LLM}}$ of 11.7 for SAE participants vs. a success rate of 39.4% with an $ECE_{\text{Human-LLM}}$ of 20.3 for AAVE participants. For AAVE participants, agents perform worse (11.2 percentage point decrease in success rate) and simulated users are more poorly calibrated (8.6 percentage point increase in ECE). In practice, such differences in performance and user simulation reliability could lead to disparities in the quality of retail assistance and the ease of interactions.

We observe contrasting age-related patterns in success rates across dialects. For SAE participants, success rates increase slightly with age (~ 3.0 percentage point increase from 18-34 to 55+ group), whereas for AAVE participants, success rates decrease with age (7.6 percentage point decrease from 18-34 to 55+ group). Notably, dialect disparities in agent performance grow larger with age: there is nearly a 12 percentage point decrease in agent performance between SAE and AAVE 35-54 groups and a 19 percentage point decrease in agent performance between SAE and AAVE 55+ groups (age-stratified GEE dialect effects: $\beta_{35-54} = 0.67$, $p = 0.01$; $\beta_{55+} = 1.24$, $p = 0.001$). The calibra-

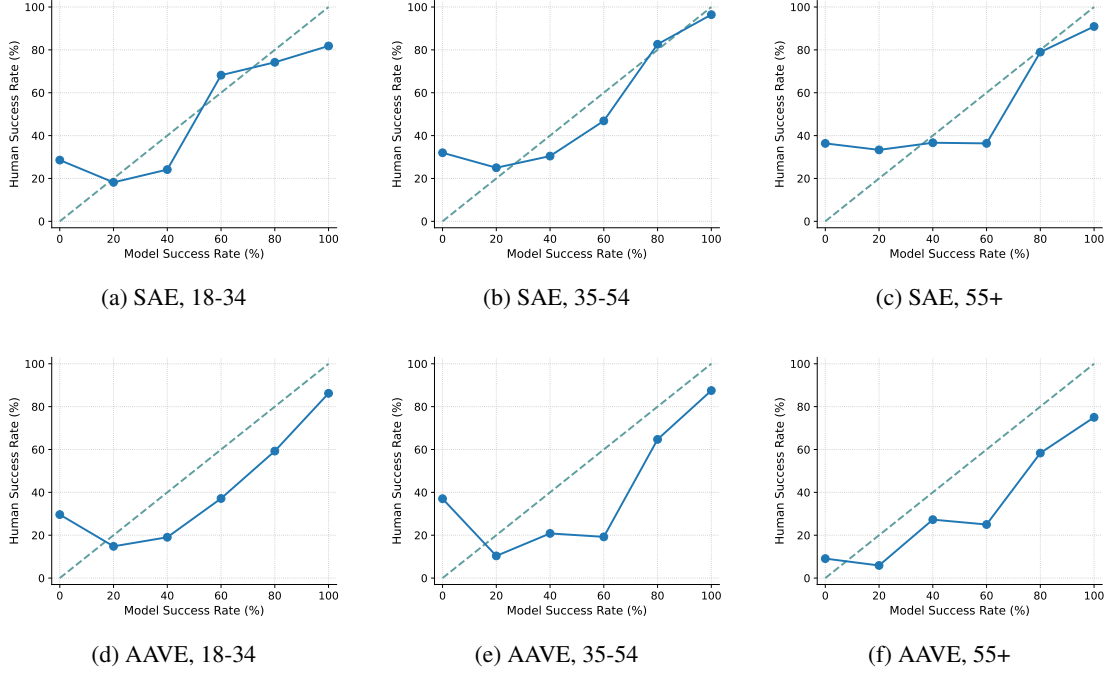


Figure 3: Human user task success vs. LLM-simulated user task success for SAE (top) and AAVE (bottom) speakers across different groups (18-34, 35-54, and 55+). The x-axis (LLM success rate) remains consistent across all panels.

tion gap is particularly pronounced for the 35-54 age group (10.3 $ECE_{\text{Human-LLM}}$ gap), where SAE primarily exhibits miscalibration for the 1st bin, while AAVE exhibits miscalibration across all bins (Figures 3b and 3e).

Note that $ECE_{\text{Human-LLM}}$ and success rate capture distinct aspects of performance. While lower $ECE_{\text{Human-LLM}}$ and higher success rate are both desirable, since lower $ECE_{\text{Human-LLM}}$ indicates that simulated users serve as reliable proxies for human users and higher success rates reflect stronger agent task performance, improvements in one do not necessarily imply improvements in the other. For example, SAE participants aged 18–34 exhibit both lower $ECE_{\text{Human-LLM}}$ and success rates, whereas SAE participants aged 55+ exhibit both higher $ECE_{\text{Human-LLM}}$ and success rates.

4.3.2 Countries

Due to participant availability constraints, we focus our cross-country analysis on the 18-34 age group. We find that differences in agent performance, shown in Table 3, are present but much less pronounced across countries (ranging from 41.0%-49.2%) than those observed by dialect and age within the US (Table 2). In particular, Kenyan and Nigerian participants experience similar success rates (43.5% and 43.7%). A GEE model confirms that cross-country differences are *not* statisti-

Group	ECE (\downarrow)	Success Rate (\uparrow)
SAE	13.0	49.2
AAVE	18.9	41.0
India	18.9	46.2
Kenya	15.6	43.5
Nigeria	17.6	43.7

Table 3: Expected Calibration Error (ECE) and success rate (%) across various dialects/countries for 18–34 age group users.

cally significant (all $p > 0.49$).

Simulated users are best calibrated to SAE participants ($ECE_{\text{Human-LLM}} = 13.0$) and worst calibrated to AAVE and Indian participants ($ECE_{\text{Human-LLM}} = 18.9$), suggesting simulated users are especially poor proxies for these groups. Across countries, we observe that simulated users tend to exhibit strongest calibration for fairly to moderately difficult tasks (20%-40% difficulty bins). However, they consistently overestimate performance for easy tasks (80%-100% difficulty bins), shown in Figure 4. As a result, evaluations relying on simulated users risk systematically underestimating difficulty agents face when deployed to diverse, global user populations.

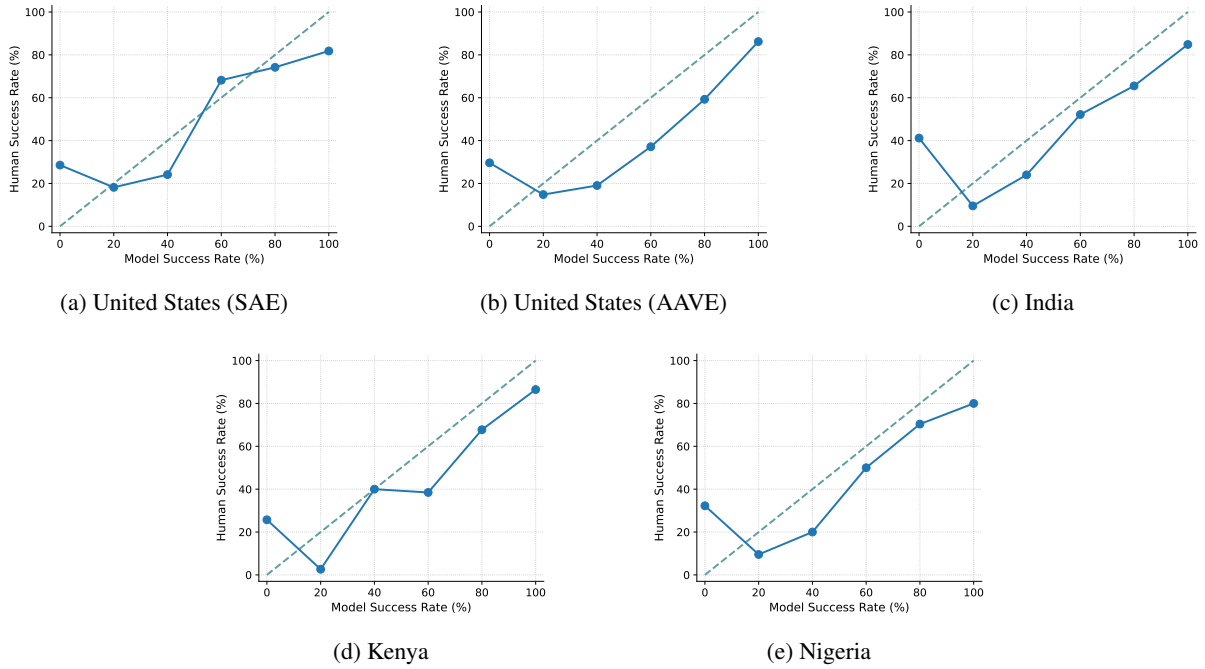


Figure 4: Human user task success vs. LLM-simulated user task success for participants across dialects/countries, all in the 18–34 age group. Note: The x-axis (LLM success rate) remains the same for all groups.

4.4 Analyzing Interactions

4.4.1 Structure and Content

Interactions between agents and simulated users vs. human users follow similar structures and surface-level forms, including # of turns, # of actions, and # of words/turn. We observe some differences between simulated and human user interactions when considering conversational content (Table 6). Simulated user conversations include questions in 18.8% of user turns and 51.8% of agent turns, compared to 9.8% and 56.3%, respectively, for human users. Notably, Nigerian participants only ask questions in 4.3% of turns.

Differences are more pronounced when examining politeness indicators (e.g., please, thank you, apologize). Simulated user conversations include such indicators in 39.2% of user and 52.0% of agent turns, compared to 19.9% and 41.1%, respectively, for human users. For both question-asking and politeness indicators, differences in behavior are more pronounced between simulated and human users than among human users. We also show that targeted behavioral interventions, such as prompting simulated users to limit politeness, can alter calibration patterns and reduce gaps for highly miscalibrated groups (see Appendix A.6), indicating that prompting strategies can partially mitigate miscalibration issues.

4.4.2 Errors

We also analyze differences in (1) error types and (2) error attribution for simulated vs. human users to understand whether both user groups experience the same failure modes. Error types include argument errors, missing actions, extra actions, and output errors (see Table 4 for definitions).

Agents make argument errors (performing the correct action with incorrect arguments) at rates of 32.2% for simulated users compared to 23.2%–45.8% for human users, with AAVE participants experiencing the highest errors. Agents tend to less frequently omit actions (17.8% vs. 15.8%–25.3%) or include unnecessary ones (5.6% vs. 7.6%–13.3%) for simulated users compared to human users. However, output errors show the reverse pattern: agents either omit or include incorrect outputs more often for simulated users (31.4%) than for human users (12.2%–23.6%). These patterns suggest that agents exhibit different behavior when interacting with simulated vs. human users; they perform more complete and efficient action sequences but make more frequent output errors.

When comparing error attribution for simulated vs. human user conversations (Table 5), we observe clear differences in where task failures occur. In simulated conversations, agents are responsible for task failures substantially more often than for human conversations (48.9% vs. 24.5%). In contrast,

Group	Arg (%)	Miss (%)	Extra (%)	Output (%)
Simulated User				
–	32.2	17.8	5.6	31.4
Human User – US, SAE				
All	25.2	18.5	9.5	16.2
18–34	27.9	15.8	7.6	18.9
35–54	24.5	18.5	10.6	17.0
55+	23.2	21.3	10.3	12.2
Human User – US, AAVE				
All	39.0	23.0	10.7	15.5
18–34	36.8	25.3	13.3	15.5
35–54	37.8	20.4	8.7	16.4
55+	45.8	24.1	9.6	13.3
Human User – Non-US, 18–34				
India	33.3	18.2	7.9	20.4
Kenya	38.5	17.2	7.8	23.6
Nigeria	33.8	22.3	10.8	21.7

Table 4: Error breakdown (%) for simulated and human users, aggregated across tasks. Error types: **argument error** (action taken matches ground truth action but with different arguments; e.g., modifies order to the wrong address), **missing action error** (ground truth action is missing from actions taken), **extra action error** (actions taken go beyond ground truth actions), and **output error** (expected outputs are missing/incorrect). Each error type is measured as a binary indicator per task (e.g., a missing action error records whether **any** required action is missing). Note that error percentages do not sum to 100%; some tasks have no errors while others have multiple error types.

Source	Simulated	Human
Agent	48.9	24.5
User	40.0	62.2
Both	2.2	11.1
Other	8.9	2.2

Table 5: Error Attribution (%) for simulated vs. human user conversations with $reward = 0$, manually annotated by the authors ($n = 45$ per condition, matched on task difficulty).

users are the primary source of failure in human conversations (62.2% vs. 40%).

These differences indicate that simulated and human interactions surface distinct failure patterns. The higher user error rate in human interactions reflects ambiguity, misunderstandings, or partial compliance that humans naturally introduce. Conversely, the higher agent error rate in simulated conversations suggests simulated users appear to exhibit more precise instruction following or adapt more readily to agent responses, placing greater burden on agents to execute correctly. In practice, this divergence may lead to misdiagnoses of failures. Simulation-based evaluations may overemphasize agent execution errors, while obscuring challenges that arise when real users engage with agents in ways not reflected by simulated users.

5 Discussion and Conclusion

As AI agents become integrated into everyday tasks, ensuring equitable performance across diverse populations is essential. However, the “user” component of agentic evaluations, central to real-world interaction, has largely been overlooked. Our findings reveal fundamental limitations of LLM-simulated users, showing that current simulation practices misestimate agent performance for actual users and obscure demographic disparities, which may result in evaluations optimizing for objectives that diverge from real-world use.

Systems optimized for simulated users may appear robust in benchmarks while failing disproportionately for real users whose communication styles are underrepresented by simulation. For example, the behavioral artifacts we observe in simulated interactions—such as heightened politeness and question-asking—suggest that simulators reflect communication norms that may not generalize across diverse user groups. Moving forward, agentic benchmarks should assess robustness across multiple simulators, validate simulated outcomes against demographically diverse human data if possible, and transparently acknowledge the limitations of user simulation.

550 Limitations

551 While our study provides important insights into
552 user simulation in agentic evaluations, it is useful
553 to clarify its scope and outline directions for future
554 work. Our evaluation is conducted entirely in En-
555 glish (replicating a limitation of popular agentic
556 benchmarks, which are limited to English), which
557 restricts our ability to assess how LLM-simulated
558 users behave in multilingual settings. Since lan-
559 guage influences user and agent behavior, capa-
560 bilities, and interaction norms, it is important to
561 verify the extent to which our findings hold for
562 non-English contexts. In addition, our age-based
563 analyses are limited to users in the United States
564 due to recruitment constraints, leaving open the
565 question of how performance and calibration dis-
566 parities vary with age across other countries and
567 cultural contexts.

568 We also evaluate agents in a single domain, fo-
569 cusing on retail customer service scenarios from
570 τ -Bench. While these tasks are designed to cap-
571 ture multi-turn, task-oriented interactions with tool
572 use, agent performance patterns and the quality of
573 user simulation may differ in other domains such
574 as healthcare, where interaction structure and task
575 complexity vary. In such domains with longer-form
576 interactions, individual style and cultural differ-
577 ences may be more apparent, potentially amplify-
578 ing the effects we observe. Nevertheless, it would
579 be important to empirically validate this expecta-
580 tion.

581 Finally, we focus on a single, fixed agent (GPT-
582 4o). Holding the agent constant enables us to iso-
583 late the effects of user simulation on evaluation
584 outcomes. However, we do not examine how cali-
585 bration gaps or performance disparities vary across
586 different agents, which is important for developing
587 a more complete understanding of the robustness,
588 validity, and fairness of user simulation. We discuss
589 these points further in Appendix A.1.

590 Ethics Statement

591 All participants received standardized compensa-
592 tion that was not adjusted by country, ensuring
593 consistent and fair payment regardless of geo-
594 graphic location. Participants were provided with
595 an overview of the study procedures upfront and
596 could withdraw from the study at any point without
597 penalty. Participants provided informed consent
598 by reading the study instructions and choosing to
599 participate and complete the study. The study is

classified as minimal risk, since it involves inter- 600
action with AI agents in simulated retail customer 601
service scenarios and does not involve the collec- 602
tion of sensitive personal information. 603

Beyond the user study itself, our work has 604
broader societal implications. Our results demon- 605
strate that LLM-simulated users may not serve as 606
reliable proxies for human users and can exhibit 607
demographic biases. If simulated users are adopted 608
as a standard practice for agent evaluation despite 609
these limitations, there is a risk that AI systems 610
could be deployed in ways that systematically un- 611
derserve certain demographic groups. We hope 612
this work encourages the research community to re- 613
flect on current evaluation practices and to develop 614
agentic evaluation approaches that better represent 615
diverse user populations. 616

References 617

- Dhruv Agarwal, Mor Naaman, and Aditya Vashistha. 618
2025. [Ai suggestions homogenize writing toward](#) 619
[western styles and diminish cultural nuances](#). CHI 620
'25, New York, NY, USA. Association for Computing 621
Machinery. 622
- Victor Barres, Honghua Dong, Soham Ray, Xujie Si, 623
and Karthik Narasimhan. 2025. [\$\tau^2\$ -bench: Evaluat-](#) 624
[ing conversational agents in a dual-control environ-](#) 625
[ment](#). *Preprint*, arXiv:2506.07982. 626
- Elisa Bassignana, Amanda Cercas Curry, and Dirk Hovy. 627
2025. [The AI gap: How socioeconomic status affects](#) 628
[language technology interactions](#). In *Proceedings* 629
of the 63rd Annual Meeting of the Association for 630
Computational Linguistics (Volume 1: Long Papers), 631
pages 18647–18664, Vienna, Austria. Association 632
for Computational Linguistics. 633
- Samuel Cahyawijaya, Delong Chen, Yejin Bang, Leila 634
Khalatbari, Bryan Wilie, Ziwei Ji, Etsuko Ishii, and 635
Pascale Fung. 2025. [High-dimension human value](#) 636
[representation in large language models](#). In *Pro-* 637
ceedings of the 2025 Conference of the Nations of 638
the Americas Chapter of the Association for Com- 639
putational Linguistics: Human Language Technolo- 640
gies (Volume 1: Long Papers), pages 5303–5330, 641
Albuquerque, New Mexico. Association for Compu- 642
tational Linguistics. 643
- Serina Chang, Ashton Anderson, and Jake M. Hof- 644
man. 2025. [ChatBench: From static benchmarks](#) 645
[to human-AI evaluation](#). In *Proceedings of the 63rd* 646
Annual Meeting of the Association for Computational 647
Linguistics (Volume 1: Long Papers), pages 26009– 648
26038, Vienna, Austria. Association for Computa- 649
tional Linguistics. 650
- CNBC. 2025. [Ai travel agents planning future trip far](#) 651
[beyond assistant status](#). 652

769	Olawale Salaudeen, Anka Reuel, Ahmed Ahmed, Suhana Bedi, Zachary Robertson, Sudharsan Sundar, Ben Domingue, Angelina Wang, and Sanmi Koyejo. 2025. Measurement to meaning: A validity-centered framework for ai evaluation . <i>Preprint</i> , arXiv:2505.10573.	825
770		826
771		827
772		828
773		829
774		830
775	Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoos Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? In <i>Proceedings of the 40th International Conference on Machine Learning, ICML'23</i> . JMLR.org.	831
776		832
777		833
778		834
779		835
780	Sebastin Santy, Jenny Liang, Ronan Le Bras, Katharina Reinecke, and Maarten Sap. 2023. NLPositionality: Characterizing design biases of datasets and models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9080–9102, Toronto, Canada. Association for Computational Linguistics.	836
781		837
782		838
783		839
784		840
785		841
786		842
787	Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The risk of racial bias in hate speech detection . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1668–1678, Florence, Italy. Association for Computational Linguistics.	843
788		844
789		845
790		846
791		847
792		848
793	Maarten Sap, Swabha Swayamdipta, Laura Vianna, Xuhui Zhou, Yejin Choi, and Noah A. Smith. 2022. Annotators with attitudes: How annotator beliefs and identities bias toxic language detection . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5884–5906, Seattle, United States. Association for Computational Linguistics.	849
794		850
795		851
796		852
797		853
798		854
799		855
800		856
801		857
802	Yijia Shao, Vinay Samuel, Yucheng Jiang, John Yang, and Diyi Yang. 2025. Collaborative gym: A framework for enabling and evaluating human-agent collaboration . <i>Preprint</i> , arXiv:2412.15701.	858
803		859
804		860
805		861
806	Huaman Sun, Jiaxin Pei, Minje Choi, and David Jurgens. 2025. Sociodemographic prompting is not yet an effective approach for simulating subjective judgments with LLMs . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 845–854, Albuquerque, New Mexico. Association for Computational Linguistics.	862
807		863
808		864
809		865
810		866
811		867
812		868
813		869
814		870
815	Yan Tao, Olga Viberg, Ryan S Baker, and René F Kizilcec. 2024. Cultural bias and cultural alignment of large language models . <i>PNAS Nexus</i> , 3(9).	871
816		872
817		873
818	Kimberly Truong, Riccardo Fogliato, Hoda Heidari, and Steven Wu. 2025. Persona-augmented benchmarking: Evaluating LLMs across diverse writing styles . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 22687–22720, Suzhou, China. Association for Computational Linguistics.	874
819		875
820		876
821		877
822		878
823		879
824		880
	Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024. Two tales of persona in LLMs: A survey of role-playing and personalization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 16612–16631, Miami, Florida, USA. Association for Computational Linguistics.	825
		826
		827
		828
		829
		830
		831
	Haoxin Wang, Xianhan Peng, Huang Cheng, Yizhe Huang, Ming Gong, Chenghan Yang, Yang Liu, and Jiang Lin. 2025a. ECom-bench: Can LLM agent resolve real-world E-commerce customer support issues? In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 276–284, Suzhou (China). Association for Computational Linguistics.	832
		833
		834
		835
		836
		837
		838
		839
	Wenxuan Wang, Wenxiang Jiao, Jingyuan Huang, Ruyi Dai, Jen-tse Huang, Zhaopeng Tu, and Michael Lyu. 2024. Not all countries celebrate thanksgiving: On the cultural dominance in large language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6349–6384, Bangkok, Thailand. Association for Computational Linguistics.	840
		841
		842
		843
		844
		845
		846
		847
	Zhefan Wang, Ning Geng, Zhiqiang Guo, Weizhi Ma, and Min Zhang. 2025b. Human vs. agent in task-oriented conversations . In <i>Proceedings of the 2025 Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region, SIGIR-AP 2025</i> , page 133–142, New York, NY, USA. Association for Computing Machinery.	848
		849
		850
		851
		852
		853
		854
		855
	Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024. Travelplanner: a benchmark for real-world planning with language agents . ICML'24. JMLR.org.	856
		857
		858
		859
	Jane Xing, Tianyi Niu, and Shashank Srivastava. 2025. Chameleon LLMs: User personas influence chatbot personality shifts . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 17314–17332, Suzhou, China. Association for Computational Linguistics.	860
		861
		862
		863
		864
		865
	Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Zhiruo Wang, Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Martin, Zhe Su, Leander Melroy Maben, Raj Mehta, Wayne Chi, Lawrence Kunho Jang, Yiqing Xie, and 2 others. 2025. Theagent-company: Benchmarking LLM agents on consequential real world tasks . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	866
		867
		868
		869
		870
		871
		872
		873
		874
		875
	Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2025. τ-bench: A benchmark for Tool-Agent-User interaction in real-world domains . In <i>International Conference on Representation Learning</i> , volume 2025, pages 9965–10017.	876
		877
		878
		879
		880

Group	Turns	Actions	W/T (U)	W/T (A)	Q (U)	Q (A)	P (U)	P (A)
Simulated User								
–	16.2	7.9	13.4	55.0	18.8	51.8	39.2	52.0
Human User – US, SAE								
All	15.0	7.4	12.6	53.0	11.6	56.6	19.1	41.1
18–34	15.8	7.6	11.4	55.4	14.5	58.1	15.6	41.1
35–54	14.9	7.3	11.6	51.5	9.8	55.3	19.7	43.6
55+	14.3	7.3	14.8	51.9	10.3	56.2	22.1	38.7
Human User – US, AAVE								
All	14.9	7.4	12.4	52.8	10.0	56.4	19.7	42.5
18–34	15.5	7.6	11.9	52.6	12.8	56.8	17.4	44.9
35–54	14.8	7.3	12.1	53.4	8.9	57.2	19.9	42.1
55+	15.3	7.0	13.9	51.9	9.1	53.9	24.0	38.3
Human User – Non-US, 18–34								
India	16.3	7.4	11.5	53.1	9.9	57.7	20.1	40.8
Kenya	14.4	8.0	13.7	54.6	8.9	54.7	21.4	39.1
Nigeria	14.4	7.5	11.2	55.6	4.3	56.8	18.7	41.4

Table 6: Conversational statistics (means) for simulated users and human users (split by demographic group). Statistics include: **Turns** – number of turns in the interaction, **Actions** – total number of actions performed by the agent (including read and write actions), **W/T (U/A)** – words per turn, for the user/agent, **Q (U/A)** – percent of turns with a question, for the user/agent, **P (U/A)** – percent of turns with politeness indicators (e.g., please, thank you, apologize, etc.), for the user/agent.

Age Group	ECE Δ
US, SAE	
All	2.4
18–34	-5.0
35–54	3.3
55+	5.5
US, AAVE	
All	-4.9
18–34	-2.6
35–54	-6.0
55+	-0.3
Non-US, 18–34	
India	-6.2
Kenya	-4.2
Nigeria	1.0

Table 7: Differences between $ECE_{\text{Human-LLM}}$ for reduced-politeness user simulation across demographic groups vs. standard user simulation. Negative values indicate better calibration after behavioral intervention. Human user study results are held fixed, while task difficulty bins are recomputed based on updated simulated user results, resulting in changes to $ECE_{\text{Human-LLM}}$.

behavioral cues (e.g., "You are detail-oriented and want to make sure everything is addressed in one go") to avoid biasing user behavior (Tseng et al., 2024; Xing et al., 2025). These modifications preserve all task objectives and requirements.

Example of Original Task Instructions *You are Yusuf Rossi in 19122. You received your order #W2378156 and wish to exchange the mechanical keyboard for a similar one but with clicky switches and the smart thermostat for one compatible with Google Home instead of Apple HomeKit. If there is no keyboard that is clicky, RGB backlight, full size, you'd rather only exchange the thermostat. You are detail-oriented and want to make sure everything is addressed in one go.*

Example of Adapted Task Instructions *You are User b63 in 19122. You received your order #W2378156 and wish to exchange the mechanical keyboard for a similar one but with clicky switches and the smart thermostat for one compatible with Google Home instead of Apple HomeKit. If there is no keyboard that is clicky, RGB backlight, full size, you would rather exchange only the thermostat. You want to make sure everything is addressed in one go. To start the conversation, say 'Hello, my email is user.b63@example.com.'*

A.5 User Study Interface and Instructions

Users interact with a Streamlit app (Figure 5) to converse with the agent and complete tasks. We provide participants with the following general instructions, which apply to all tasks in addition to the task-specific instructions shown in Table 8.

Instructions: Please respond as the user described in the task instructions. You want to complete all the requests mentioned in the instructions. The agent is there to assist you with completing the task. Do not make up information beyond what the instructions provide. You can tell the agent you are unsure and ask them to look up information based on your profile or orders. Beyond this, please behave naturally and converse as you normally would. Use the ‘End Conversation’ button in the left sidebar to finish your conversation. To begin the conversation, authenticate yourself by providing the user email provided in the instructions.

Note that all task instructions are free of user persona information and use generic user IDs to avoid biasing interactions.

A.6 Behaviorally Targeted User Simulation

We identify heightened politeness as a behavioral artifact in LLM-simulated user interactions and examine whether explicitly constraining this behavior affects evaluation outcomes. To test this, we introduce a minimal intervention to the user simulation prompt that limits the use of politeness indicators while preserving task objectives: “You may include politeness indicators (e.g., please, thank you, sorry) occasionally, but use them sparingly. Limit their use to at most one or two times across the entire conversation.” We re-run user simulation on the 18-task subset using GPT-4o as both the agent and user LLMs, and re-bucket tasks based on the updated performance distribution. We find that 11 of 18 tasks shift difficulty bins, indicating that task difficulty estimates are sensitive to behavioral prompting, and overall success decreases from 50.0% to 46.7%.

Comparing original and updated $ECE_{\text{Human-LLM}}$ values, we observe reduced miscalibration for AAVE, Indian, and Kenyan participants—groups that previously exhibited the largest calibration errors—while miscalibration increases for SAE and Nigerian participants (Table 7). These results highlight the sensitivity of user simulation to prompt-level choices and suggest that targeted behavioral interventions can meaningfully alter calibration patterns.

AI Agent Customer Assistance

You will chat with an AI agent that provides retail assistance. The agent is there to help you complete the task and can look up information and make changes to your orders. Please follow the task instructions carefully and converse with the agent to complete the task. Use the 'End Conversation' button in the left sidebar to finish (it will appear when you begin the task). Make sure to download your conversation logs for each task and upload them in the form. You will also answer a question about each task in the form. The question will be shown in the left sidebar after task completion, and you will answer it in the form. Make sure you answer the question before moving on to the next task, you won't be able to access it later.

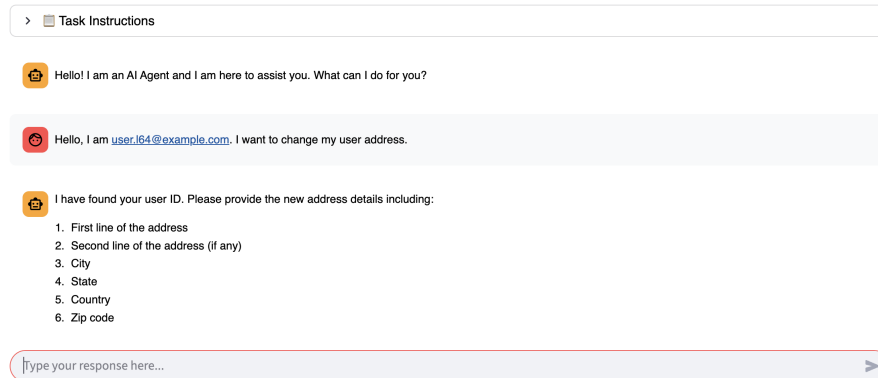


Figure 5: User study chat interface where participants interact with the GPT-4o agent to complete tasks.

Bin	Success	Example Task
1	0%	Your name is User e70 and your zip code is 32190. You just bought a water bottle with 500ml but you regret it, and you instead want to change it to the other bottle you recently ordered with 1000ml capacity. If the exact 1000ml bottle is not available any more, you can allow the material to be different. To start the conversation, say 'Hello, my email is user.e70@example.com.'
2	20%	You are User l64, and you live in Denver, 80280. You just won a lottery, and you want to upgrade all your items to the most expensive options (but make sure the shoe is still the same size). You want to pay the difference with your gift card, but if it is impossible, PayPal is fine. To start the conversation, say 'Hello, my email is user.l64@example.com.'
3	40%	Your name is User u52, and you live in 46236. Your email is user.u52@example.com. You just placed an order but you realize that your card has only \$1131 credit left, and the order total is more than \$1160. You wonder if the agent can help split the payment with another card. If not, you wonder what the most expensive item is and its price, and if you can just cancel that item. If not, you wonder if you can switch all items to their cheapest options and bring the cost down to \$1131. If so, do it. If not, you wonder if the agent can just cancel the order so that you can order again. To start the conversation, say 'Hello, my email is user.u52@example.com.'
4	60%	You are User i49, and you live in 32286. You want to exchange your skateboard for a shorter bamboo material one. If several options are available, you want to know all options and their prices, and choose the most expensive one because you believe price is quality. Also, you want to exchange the garden hose you received to the type that you just ordered (pending). To start the conversation, say 'Hello, my email is user.i49@example.com.'
5	80%	You are User b63 in 19122. You received your order #W2378156 and wish to exchange the mechanical keyboard for a similar one but with clicky switches and the smart thermostat for one compatible with Google Home instead of Apple HomeKit. If there is no keyboard that is clicky, RGB backlight, full size, you would rather exchange only the thermostat. You want to make sure everything is addressed in one go. To start the conversation, say 'Hello, my email is user.b63@example.com.'
6	100%	You are User p59, residing in Philadelphia 19031. You want to change the Desk Lamp in order #W9300146 that you've placed for the cheapest Desk Lamp that's available. Any price difference should go to a gift card. You also want to know how much you get back in total. To start the conversation, say 'Hello, my email is user.p59@example.com.'

Table 8: Example tasks from τ -Bench across difficulty bins (as measured by simulated user success across 5 runs).