ChatBench: From Static Benchmarks to Human-AI Evaluation

Anonymous ACL submission

Abstract

With the rapid adoption of LLM-based chatbots, there is a pressing need to evaluate what humans and LLMs can achieve together. However, standard benchmarks, such as MMLU, 005 measure LLM capabilities in isolation (i.e., "AIalone"). Here, we design and conduct a user study to convert MMLU questions into user-AI conversations, by seeding the user with the question and having them carry out a conversation with the LLM to answer their question. We release ChatBench, a new dataset with AI-011 alone, user-alone, and user-AI data for 396 012 questions and two LLMs, including 144,031 answers and 7,337 user-AI conversations. We find that AI-alone accuracy fails to predict user-016 AI accuracy, with significant differences across multiple subjects (math, physics, and moral 017 reasoning), and we analyze the user-AI conversations to provide insight into how they diverge 020 from AI-alone benchmarks. Finally, we show that fine-tuning a user simulator on a subset of 021 ChatBench improves its ability to estimate user-022 AI accuracies, increasing correlation on heldout questions by more than 20 points, creating possibilities for scaling interactive evaluation.

1 Introduction

027

033

037

041

In 2024, nearly 40% of US adults reported using generative AI in their everyday lives, an unprecedented rate of adoption for a new technology (Bick et al., 2024). As these models, particularly large language models (LLMs), become more integrated into our lives, it becomes increasingly important to evaluate them based on not only their capabilities in isolation, but also their interactions with humans. However, there is a large gap between human interactions and how standard benchmarks, such as Massive Multitask Language Understanding (MMLU), evaluate models (Hendrycks et al., 2021). These benchmarks test models on a fixed set of questions, and for each question, they prompt the model with the entire question text and often constrain it to respond with a single multiple choice option as its answer. In contrast, interactions with human users are far more variable, open-ended, and subject to ambiguity. Even conditioned on the same underlying intent, users may phrase their prompts differently, leave out information in their early prompts, or rely on context in later prompts. Robust AI models need to understand how to work with users in these contexts to provide accurate information and complement human expertise. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

Recently, there have been efforts to evaluate LLMs in terms of their interactions, such as evaluating real-world interactions using a strong LLM as a judge (Lin et al., 2024; Li et al., 2024c). However, these new evaluations have been largely disconnected from standard benchmarks, which are widely used; for example, every LLM released by OpenAI, Google, and Meta, inter alia, has reported its performance on MMLU (OpenAI, 2023; Gemini Team Google, 2023; Llama Team, AI@Meta, 2024). This disconnect is due to a large distribution shift between benchmark questions and questions asked by real-world users, missing the user's true intent, and missing ground-truth labels to judge the interaction, necessitating techniques like LLM-asjudge. As a result, it is difficult to directly compare results from standard benchmarks to real-world interactions or to understand how incorporating interactions changes evaluation insights.

Here, we seek to bring these lines of research closer together by directly *converting* benchmarks into user-AI conversations. We focus on MMLU, as one of the most widely used benchmarks, and design a user study where we seed users with an MMLU question and have them carry out a conversation with an LLM with the intent of answering that question. For each question, we test the LLM in isolation (i.e., "AI-alone") and evaluate the accuracy of a user interacting with the LLM (i.e., "user-AI"); furthermore, we also gather "user-alone" data per question to understand how much users improve with the LLM. This parallel data has two advantages: first, we can now conduct an applesto-apples comparison of AI-alone performance, as reported in most papers, vs. user-AI performance on the same questions, so that we can isolate the effects of incorporating interaction into evaluation. Second, recent works have explored the possibility of *simulating* the user in user-AI conversations (Li et al., 2024a) but lack sufficient data for training and testing. Our approach of "seeding" users with a question corresponds naturally to a new way to initialize user simulators, and the large-scale data we collect enables fine-tuning and validating a user simulator on this task, improving the trustworthiness of simulations for AI evaluation.

084

100

101

102

103

105

106

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

Our resulting dataset ChatBench, which we release publicly, consists of AI-alone, user-alone, and user-AI data for 396 questions and two LLMs (GPT-40 and Llama-3.1-8b), with 144,031 answers and 7,337 user-AI conversations. Our study design also includes two user-AI conditions—where the user attempts the question first on their own vs. uses AI directly—to explore nuances in user behavior. Our study reveals that AI-alone accuracy fails to predict user-AI accuracy, with significant differences across multiple subjects (math, physics, and moral reasoning). We also analyze the user-AI conversations to understand where user-AI interactions are diverging from AI-alone benchmarks. Our contributions are summarized as follows:

- We design and conduct a user study to convert MMLU questions into user-AI conversations and release a large-scale dataset ChatBench.
- We show that AI-alone accuracy fails to predict user-AI accuracy, across subjects, models, AI-alone methods, and user-AI conditions, and we analyze user-AI conversations to understand where AI-alone and user-AI diverge.
- We develop a new user simulator that mimics our user study task and show that fine-tuning our simulator on ChatBench improves its correlation with real user-AI accuracies by 21-27 points and outperforms baselines.

All together, our work helps to reconcile two vital lines of research in AI evaluation, revealing how interactions change evaluation insights and paving the way towards scalable interactive evaluation.¹

2 Related Work

Benchmarks. In this work, we focus on MMLU as one of the most commonly used LLM benchmarks (Hendrycks et al., 2021). MMLU is a question-answering (QA) dataset, consisting of multiple choice questions across 57 subjects (which we discuss in detail in Section 3.2). We also draw on the efforts of MMLU-Redux (Gema et al., 2024), where authors noted some quality concerns in the original MMLU, so they sampled a large number of MMLU questions and manually annotated them for errors. While we conduct our user study on MMLU, our approach of converting QA benchmarks to a user-AI conversation is general, and could be applied to other QA benchmarks, such as HotPotQA (Yang et al., 2018) or GSM8K (Cobbe et al., 2021), as well as adapted to non-QA tasks.

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

Evaluating human-AI interactions. Recently, there have been growing efforts to evaluate AI models based on their interactions with humans. For example, some works gather real-world interactions (e.g., WildChat (Zhao et al., 2024), ChatbotArena (Chiang et al., 2024)) and evaluate the interactions (e.g., WildBench (Lin et al., 2024), ArenaHard (Li et al., 2024c), MT-Bench (Zheng et al., 2023)), typically using a strong LLM as a judge. However, as discussed before, it is difficult to directly compare these evaluation results to standard benchmarks, due to the distribution shift in questions and change in evaluation metric. Other works have evaluated human-AI interactions in diverse contexts, such as theorem proving (Collins et al., 2024), co-writing with AI (Shen and Wu, 2023), and education (Jurenka et al., 2024), and sought to understand where human-AI combinations outperform either alone (Bansal et al., 2021; Vaccaro et al., 2024).

Our work builds on Lee et al. (2023), who make a strong argument for the need to evaluate human-LM interactions, covering five types of tasks including QA. Their work includes an exploratory user study where they have users interactively answer MMLU questions; however, they only test 30 questions and do not explore simulation. Our study greatly builds on theirs by testing 396 questions, at a large enough scale to estimate significant effects and fine-tune a user simulator, and introduces an AI-alone method that is a far more realistic proxy of a user's experience. Furthermore, our study tests more sophisticated LLMs, complex reasoning subjects, and user-AI effects across levels of question difficulty and different user-AI conditions. Our

¹Our code is available at https://anonymous.4open. science/r/interactive-eval-4813. Our dataset Chat-Bench will be made available upon publication.

work is also similar in spirit (although different in domain) to Li et al. (2024b), who convert medical
benchmarks into simulated interactions between a
patient and an expert.

Simulation with LLMs. LLMs have shown 185 promising capabilities to realistically simulate hu-186 187 man behaviors, such as responses to surveys and social science experiments (Argyle et al., 2023; Horton, 2023; Hwang et al., 2023; Hewitt et al., 2024) or interactions between humans (Park et al., 2023; Chang et al., 2024). There is also much interest 191 192 in developing LLM-based user simulators to scale AI evaluation and training (Dubois et al., 2023; 193 Ren et al., 2024; Kong et al., 2024; Li et al., 2024a). 195 However, LLMs can sometimes produce unrealistic simulations of humans, with risks of stereotyping, 196 bias, or uniformity (Cheng et al., 2023a,b; Bisbee et al., 2024; Wang et al., 2024). Thus, there is a 198 need to rigorously test whether LLM simulators produce realistic outputs and match insights that we would learn from real humans. Here, we examine a setting with well-defined simulator goals (i.e., does the simulator match user behavior and accuracy in real user-AI conversations) and release 204 a large-scale dataset that enables training and validation of simulators in this setting.

3 User Study Design

209

210

212

213

214

215

216

218

219

220

221

In this section, we discuss our user study design, including the task flow and interface, how we selected questions, and data collection. We provide additional details in Appendix A.

3.1 Task Flow and Interface

Figure 1 shows the flow of our user study. In Phase 1, users are asked to answer each question to the best of their ability on their own. In Phase 2, users are asked to chat with an unnamed "AI Chatbot" to help them answer their question. We test two LLMs, contrasting GPT-40 as a strong model and Llama-3.1-8b as a relatively weaker model. We require interaction in Phase 2-the user cannot move onto the next question without sending a message and we say that low-effort conversations, e.g., only "hi", will be flagged-but otherwise, we do not specify at all how the user should interact with the AI Chatbot. In both phases, users are asked to first report how confident they are about approaching the problem, before attempting to answer it. This additional question-level variable allows us to analyze how AI assistance helps users across varying



Figure 1: Flow of our user study.

PHASE 2 Question 3 of 9	AI Chatbot
Now, try to answer this question with the help of the Al Chatbot, using it as you would normally use any chatbot. The Al chatbot can only see the messages you send to it and not the question on the screen. You may say anything you want to the Al Chatbot to help you with answering the question. A right triangle has legs of length 6 and 8, and a hypotenuse of length c. The permeter of the triangle is 18. Compute c.	You how do i get hypotenuse length Bot To find the length of the hypotenuse of a right margie, you can use the Pythagorean theorem. The Pythagorean theorem tates that the organ. The Pythagorean theorem tates that in a right triangle, the square of the length of the hypotenuse (the side opposite the right angle) is equal to the squares of the lengths of the other two scales. The formulae
 ○ A. 18 ○ B. ¹⁵/₂ ○ C. 8 ○ D. 6 	is: $c=\sqrt{a^2+b^2}$ where: - c is the length of the hypotenuse, - a and b are the lengths of the other two sides.
CONTINUE	Here are the steps to find the hypotenuse length: 1. Square the lengths of the two sides. 2. Add the squares together. 3. Take the square root of the sum. You can send 19 more messages to Al.

Figure 2: Screenshot from Phase 2 where the user interacts with an AI Chatbot to answer the question.

230

231

232

233

234

235

236

237

238

240

241

242

243

244

245

246

247

248

249

levels of confidence. After Phase 2, all users provide feedback on the task, with free-text responses including whether they found the AI Chatbot helpful and if they saw it make any mistakes. In Figure 2, we show a screenshot of what users see in Phase 2; in the Appendix, we provide screenshots of all other pages in our task (Figures A2-A9).

Conditions. We explore two user-AI conditions: *answer-first* and *direct-to-AI*. In the *answer-first* condition, the user attempts to answer each Phase 2 question on their own first before answering with AI, but in the *direct-to-AI* condition, they have immediate access to AI for the Phase 2 questions (in both conditions, Phase 1 is all user-alone). The advantage of *answer-first* is that, for the same question, we can record a user's answer on their own vs. with AI, allowing us to estimate the marginal impact of AI more precisely (i.e., within-subjects), while for *direct-to-AI*, the set of user-alone answers and user-AI answers for a given question come



Figure 3: Examples of questions from our user study.

from different users (i.e., between-subjects). However, we hypothesized that user behavior and accuracy in the user-AI stage could be impacted by the user attempting the answer first, reducing ecological validity if we believe users typically go directly to AI in the real world. Thus, we keep both conditions, allowing us to test our hypothesis and explore nuances in user behavior.

250

251

258

259

261

263

264

268

269

271

272

275

277

278

Incentivization. To incentivize participants in our study to answer questions correctly, we included a small bonus of \$0.10 per correct answer, on top of a base compensation of \$5.00 for completing the entire task. We included these incentives to improve ecological validity, since our study is meant to capture how a real-world user would interact with an AI system if they have a question that they genuinely want to answer. In Appendix A.1, we discuss pilots we ran with and without incentivization, as well as how we mitigated risks of cheating with external tools.

3.2 Question Selection

We consider five datasets from MMLU for our experiments: Elementary, High School, and College Mathematics, Conceptual Physics, and Moral Scenarios. We include three math datasets since this subject still poses unique challenges for LLMs: for example, the HELM leaderboard (Liang et al., 2023) reports that while GPT-4o's mean accuracy on MMLU is 84%, its accuracy is only 48% on High School Math and 51% on College Math.²

Furthermore, the three math datasets stratify different levels of difficulty for humans, allowing us to explore how user-AI effects change across difficulty levels. We also include Conceptual Physics and Moral Scenarios as two other reasoning domains with very different types of problems and differing levels of human expertise. In Figure 3, we provide examples of questions from these datasets, showcasing their diversity. 280

281

282

285

287

289

290

291

292

293

294

295

296

297

298

299

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

To aid with question selection, we use the annotations from MMLU-Redux (Gema et al., 2024). The authors recognized occasional quality issues with the original MMLU, so for each MMLU dataset, they sampled 100 questions from the test set uniformly at random and labeled them for errors. While they found many errors in some datasets (e.g., Virology), the majority of the questions (92%-99%) in the datasets we chose passed their review. As a second layer of quality control, we also ran OpenAI's advanced reasoning o1 model over the 100 questions and manually checked the questions that o1 did not get correct. We kept the intersection of questions that passed MMLU-Redux's inspection and ours (with o1's help).

To reduce variance in the number of Batches. answers that each question received, we organized the questions into *batches* and selected a random batch per user, instead of selecting each question randomly. For the math questions, each batch consisted of 5 elementary, 5 high school, and 2 college questions. We included fewer college questions since we found in pilots that college questions were too difficult for most users, so they tended to defer to the LLM's first answer without much interaction. Based on the number of questions that passed inspection, we were able to create 19 math batches, with 95 elementary, 95 high school, and 38 college questions in total. For Conceptual Physics and Moral Scenarios, we constructed 7 batches of size 12, resulting in 84 questions for each subject.

3.3 Data Collection

We recruited workers on Prolific to participate as users in our study (see eligibility criteria in Appendix A). For our full pre-registered study, we recruited 650 workers, and we also ran two mediumsized pilots (100 workers without incentives and 60 workers with incentives). When a user began the study, they were randomly assigned to one of the three subjects (60% probability for math, 20% for conceptual physics, and 20% for moral scenarios)

²https://crfm.stanford.edu/helm/mmlu/latest/#/ leaderboard

428

378

379

380

381

382

and assigned uniformly at random to one of that subject's question batches, one of the two user-AI conditions, and one of the two models (GPT-40 and Llama-3.1-8b). Within the question batch, 3 questions were randomly assigned to Phase 1 and 9 to Phase 2. We also included an attention-check question for every user, which we found the vast majority (over 99%) of users passed.

Compiling data over the three runs, we have 10,831 confidence answers, 7,143 user-alone answers, and 7,337 user-AI answers and conversations in ChatBench (see Table A3 for additional data statistics). While we include data from all three runs in ChatBench to provide a larger resource for the community, for our analyses in the rest of the paper, we only use data from the workers in our full pre-registered study so that populations within our analysis are entirely comparable.

4 **Experimental Results**

335

336

341

344

347

349

357

370

372

374

375

376

377

In this section, we describe our experimental results, including how we conducted AI-alone experiments, comparisons of AI-alone vs. user-AI results, and analyses of the user-AI conversations. For our main results comparing AI-alone vs. user-AI, we preregistered our analyses on AsPredicted.³ We provide additional results and methodological details (e.g., statistical tests) in Appendix B.

AI-Alone Experiments 4.1

Our goal in this work is to understand how evaluation conclusions change when we move from AI-alone to user-AI settings. However, even for a fixed benchmark, there can be multiple ways to evaluate an LLM on its own. First, we try letteronly methods, which require the model to answer with only a single letter corresponding to the selected answer option ("A" through "D"). This is the method used by Lee et al. (2023), along with various leaderboards, such as HELM (Liang et al., 2023), to standardize the answer format. We try two letter-only variants, zero-shot and few-shot, where we prepend the 5 examples from the MMLU "dev" set to the prompt as in-context examples.

We also introduce a more realistic AI-alone technique which serves as a better proxy for user experience by not constraining the model's response format. The method, which we call *free-text*, is very simple: (1) prompt the evaluated model with the concatenated question text and answer options,

solute deviation of 21 percentage points, averaged over the 10 dataset and model pairs. With a few exceptions-specifically Conceptual Physics for Llama-3.1-8b and College Mathematics and Moral Scenarios for GPT-40-all differences are statistically significant. Results are similar for zero-shot letter-only, which we report in Tables A1-A2. Notably, our AI-alone method, free-text (dark blue), is a much better predictor of user-AI accuracy, reducing the mean absolute deviation to 10 percentage points. However, it still differs significantly from user-AI performance, notably for Moral Scenarios with Llama-3.1-8b and for all datasets except

Our results also reveal the complexity of combining humans and AI, as the size of gaps and ordering between user-alone, user-AI, and AI-alone vary over models and datasets. For example, for the math datasets, GPT-40 performs quite well on its own (using free-text), while humans struggle on their own, especially for high school and college. In these cases, user-AI accuracy is between the two, significantly better than user-alone and significantly worse than AI-alone. Meanwhile, Llama-3.1-8b performs significantly worse than GPT-40 on the math datasets, but we do not see a further drop in performance from AI-alone to user-AI. In the fol-

without any additional instructions, (2) use GPT-40 to extract an answer (if any) from the response. We include the full prompts for all three AI-alone methods in Listings 1-4.

We ran these three AI-alone methods on the two models and all 396 questions from our user study, gathering 50 answers per model and question. As shown in Figure 4, our few-shot letter-only results for GPT-40 approximately match those reported on the HELM leaderboard per dataset (which is also few-shot letter-only, but uses the entire test set). While prior work, like HELM, often uses temperatures of 0 for multiple choice QA, we used a temperature of 0.7, since we wanted to perfectly match the model parameters used in the user study, and 0.7 is a more realistic temperature for realworld AI chatbots.

4.2 AI-Alone vs. User-AI

Dataset-level accuracy. We visualize our main results in Figure 4, which shows mean accuracy per model and dataset, over user-alone (red), user-AI (purple), and AI-alone (blue). First, we see that few-shot letter-only (light blue) is a very poor predictor of user-AI performance, with a mean ab-Moral Scenarios with GPT-40.

³https://aspredicted.org/n84n-sn3f.pdf.



Figure 4: Mean accuracy per model and dataset, comparing user-alone (red), user-AI (purple), AI-alone free-text (dark blue), and AI-alone letter-only few-shot (light blue). See Tables A1-A2 for numbers and statistical tests.

lowing section, we uncover counterveilling factors 429 that explain these results: on one hand, users in-430 troduce ambiguity compared to AI-alone methods, 431 which include the entire question text and answer 432 options; on the other hand, users can sometimes 433 recognize mistakes in AI reasoning, of which there 434 are more for Llama-3.1-8b. Finally, our results re-435 veal that even when AI-alone benchmarks report 436 437 a large gap in performance between two models, this gap can become much smaller after incorpo-438 rating user interactions. Comparing GPT-40 and 439 Llama-3.1-8b, their average gap in AI-alone free-440 text accuracy is 25 percentage points, but this gap 441 shrinks to less than 10 percentage points in user-AI 442 interactions (9 percentage points for direct-to-AI 443 and 5 percentage points for answer-first). 444

Question-level accuracy. Besides mean accu-445 racy, we can also measure the correlation in per-446 question accuracies. We find that the Pearson cor-447 relation between AI-alone free-text and user-AI 448 is only r = 0.45 for *direct-to-AI* and r = 0.46449 for answer-first. While correlations may be lower 450 because per-question user-AI accuracies are imper-451 fectly measured, the free-text correlation is still 452 well below what we would expect if user-AI ac-453 curacies were drawn from the same distribution 454 as free-text, which would range from r = 0.88 to 455 0.94 (Section B.2). We also examine the correla-456 457 tion with per-question differences in user-AI and user-alone accuracy, since it may be more reason-458 able to expect AI-alone to predict the improvement 459 the user makes with AI assistance, instead of the 460 overall accuracy. However, the correlations remain 461



Figure 5: Fraction of user-AI interactions that mirror AI benchmark, by subject and model.

low, at r = 0.26 for *direct-to-AI* and r = 0.27for *answer-first*, suggesting that AI-alone cannot predict improvement very well either. Similarly, user-AI accuracy cannot be reliably predicted from user-alone and AI-alone accuracies at the question level. A linear model yields a correlation of 0.55 for predicting user-AI accuracies from free-text and *answer-first* accuracies, and 0.63 when using free-text and *direct-to-AI* accuracies.

462

463

464

465

466

467

468

469

470

471

4.3 Characterizing User-AI Conversations

Our summary results show that user-AI accuracies472are significantly different from AI-alone accuracies.473To better understand what drives these differences474we use a separate LLM as an annotator to charac-475terize the user-AI conversations. For each user-AI476conversation, we gather the full log of the conversa-477tion and its associated metadata (e.g., the question478

ID, the correct answer, the user's selected answer, etc.), and prompt a separate instance of GPT-40 to use this information to extract the answers to several classification questions: whether the first substantive user prompt was a question, whether the first user question was a near-exact rephrasing of the original question or one of several other possibilities, and whether the first and last AI answers were correct (Listing 5).

479

480

481

482

483

484

485

486

487

488

489

491

492

493

494

495

496

497

498

499

502

503

504

506

511

512

513

514

515

516

517

518

519

521

523

525

527

How often does the conversation follow what we might expect if AI benchmarks were faithful proxies of human-AI interaction? We say a conversation *mirrors* an AI benchmark if (1) the user's first substantive prompt is a near-exact rephrasing of the question (otherwise the user is injecting their own knowledge or information into the interaction), (2) the LLM responds with an answer, and (3) the user submits that answer. In Figure 5, we see that only 34% of all interactions mirror AI benchmarks, revealing the extent to which user-AI interactions diverge from AI benchmarks. Among the remaining interactions, we find that a primary source of divergence is the user asking a related but different question, which is often ambiguous (e.g., leaving out critical information for a math problem). On the other hand, we find that users occasionally correct the AI model's wrong answers, especially with the weaker model, Llama-3.1-8b (Figure B1).

Using data from the *answer-first* condition also reveals that AI helps humans more often than it hinders them. When the same user answers a question first without AI and then with AI assistance, more than half (54%) of incorrect user-alone answers are corrected with AI support, while only 10% of correct user-alone answers turn incorrect with AI assistance.

5 Simulating User-AI Conversations

From our user study, we showed that incorporating user interactions significantly changes evaluation conclusions, compared to AI-alone evaluation. However, data from human users is costly and timeconsuming to collect, motivating the development of a user simulator to scale interactive evaluation. In this section, we describe our user simulator and present experimental results.

5.1 Fine-Tuning a User Simulator

We define a new user simulator that we can finetune on our collected user data, by mimicking the experience of users in our study. First, we seed



Figure 6: Example of prompts to our two-step user simulator, using one of the example questions from Figure 3. See Listings 6-8 for complete prompts.

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

the user simulator with the MMLU question, as we did with human users in our study, and we tell the simulator to interact with an AI system to answer its question (Figure 6, left). Then, we break the simulator's task into two subtasks: (1) when there is no conversation yet, we prompt the simulator to generate its first prompt as a user (Figure 6, top right), (2) given the conversation so far, we prompt the simulator to either answer the question in the form "Answer: LETTER", if the question has been answered by the conversation, or if not, generate the next user prompt (Figure 6, bottom right).

We then transform the real user-AI conversations from our study into training examples for supervised fine-tuning. Each conversation with kuser utterances yields k + 1 training examples: one example in the Task 1 format where the gold standard response is the real user's first utterance; k - 1examples in the Task 2 format where the gold standard response is each of the remaining utterances (providing the conversation up to that utterance); and one example in the Task 2 format with the full conversation and the gold standard response being "Answer: LETTER" corresponding to the user's selected multiple choice option.

5.2 User Simulator Experiments

For these experiments, we use GPT-40 as our simulator. We try four baselines: the two AI-alone methods, the two-step simulator without fine-tuning, and the user simulator from IQA-EVAL (Li et al., 2024a). Their simulator, designed with prompt engineering, receives a prompt consisting of a role description ("You are mimicking a human."), a task

			AI:	GPT-40			AI: Lla	ama-3.1-8	b
Type	Method	Corr. ↑	$MAE\downarrow$	BLEU ↑	ROUGE ↑	Corr.	MAE	BLEU	ROUGE
AI-alone	Letter-only few-shot	0.33	0.31	_	-	0.24	0.41	-	—
AI-alone	Free-text	0.44	0.22	_	-	0.57	0.22	_	_
Sim-AI	IQA-EVAL	0.43	0.20	0.085	0.311	0.51	0.19	0.086	0.313
Sim-AI	Two-Step	0.36	0.22	0.102	0.347	0.42	0.21	0.102	0.346
Sim-AI	ChatBench-Sim	0.63	0.15	0.261	0.460	0.63	0.17	0.258	0.457

Table 1: Comparing to user-AI conversations: AI-alone methods, IQA-EVAL (Li et al., 2024a), and the two-step simulator before (Two-Step) and after fine-tuning on ChatBench (ChatBench-Sim). Top-performing is bolded.

description ("You are trying to choose the correct answer for the given question."), and discussion instructions (e.g., "In each turn, please only ask one sub-question to interact with the assistant."); see Listing 9 for the full prompt. We compare these baselines to our model, the two-step simulator finetuned on ChatBench ("ChatBench-Sim").

561

562

563

565

570

571

572

574

575

576

577

578

In our fine-tuning experiments, we randomly split the questions from our user study into 60% for training (n = 237) and withheld 40% for testing (n = 159), and we fine-tuned on all user-AI conversations for the train questions. For all three simulator methods, we test them on the held-out test questions by generating conversations entirely from scratch, given only the question (in contrast, an easier but less realistic set-up would be to provide the real conversation up to the n^{th} turn and have the simulator generate the next user utterance).

Evaluation metrics. We generate 10 simulator-AI conversations per test question and compare to real user-AI conversations for the same question and AI system. To evaluate whether accuracies are similar, we measure the correlation and mean abso-583 584 lute error (MAE) between simulator-AI vs. user-AI accuracies, only keeping test questions where we have at least 10 user-AI answers (n = 146). To 586 evaluate whether the simulator's generated utterances are realistic, we measure the average BLEU and ROUGE scores of the simulator's first prompt 589 compared to the real user's first prompt. 590

591**Results.** As shown in Table 1, fine-tuning our592simulator yields large gains, with a 22-27 point593increase in correlation and a 30-52% decrease in594MSE. As shown in Figures B3-B4, a primary fail-595ure mode of the simulator before fine-tuning is that596it cannot replicate human mistakes and greatly over-597estimates user-AI performance, producing far more598questions with accuracies of 1.0 than we see in the599real user-AI distribution, while the fine-tuned sim-600ulator matches the real distribution more closely.601We also find that fine-tuning improves the real-

ism of the simulator's generated utterances, with 11-16 point improvements in BLEU and ROUGE. The fine-tuned simulator also outperforms both AIalone methods and IQA-EVAL across metrics. 602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

6 Conclusion

We have shown that evaluation conclusions change significantly from AI-alone benchmarks to user-AI interactions, across question domains, AI models, AI-alone methods, and user-AI conditions. Our results motivate the need for more realistic evaluations of AI models that incorporate user interactions. However, this goal is difficult to achieve, as user data is expensive to collect. To make this goal more feasible, we both release a new largescale dataset of user interactions, ChatBench, and demonstrate the potential of building user simulators to scale interactive evaluation.

The changes we see from AI-alone to user-AI accuracies are often large enough to affect qualitative conclusions about the models. For example, what can seem like a large disparity between models on AI-alone benchmarks (e.g., 25 percentage point gap between GPT-40 and Llama-3.1-8b on free-text) can shrink to much smaller gaps after incorporating user interactions (e.g., 5 point gap for answer-first). These changes could impact realworld decisions, such as which model to deploy (e.g., a lightweight, on-device model that performs only slightly worse than a much larger off-device model might be preferable in some circumstances). To this end, in future work we hope to understand how AI-alone benchmarks are currently used to make decisions and how those decisions might change after taking into account human interactions. We also hope to expand our analysis to more benchmarks and non-QA tasks. Finally, we hope to develop training techniques to build even more realistic user simulators: while we see large improvements from fine-tuning on ChatBench, the best correlations only reach 0.63, leaving room for future improvement and innovation.

7 Limitations

643

688

Our work has several limitations, which we tried to mitigate but should be taken into consideration when interpreting the results.

Coverage. Our user study has limited coverage 647 of possible benchmarks and user tasks. We chose to focus on the MMLU benchmark (Hendrycks et al., 2021) and question-answering as our task, since 651 MMLU is one of the most popular LLM benchmarks and it covers a wide range of subjects, so we could test multiple subjects in comparable ways and with minimal changes to our user study. We began with question-answering since we can naturally transform a benchmark question into a user-AI conversation, where the user is trying to answer the question. However, future work should investigate whether results are consistent on other benchmarks and/or tasks, especially more open-ended generation tasks that are common in real-world user-AI interactions (Zhao et al., 2024).

Ecological validity. Our user study is meant to capture how a user would act if they have a question in mind and they are interacting with an AI system to answer their question. However, since we wanted to match the user's underlying question with the MMLU questions, we had to tell the user what question to answer, which could lead to different behavior compared to if they were intrinsically motivated to answer a question. To mitigate this, 671 we included a small incentive (\$0.10 per correct answer), so that they would try to get the correct 673 answer, and we filtered out users who failed the at-674 tention check; however, it is still possible that users' 675 behaviors would be different in the real world. Our study setting was also different from real world question-answering: we recruited workers on Prolific to do our study, where they answered 13 ques-679 tions consecutively in our interface. Still, we tried to match real-world settings, such as choosing models they might interact with in the real world (e.g., GPT-40), using realistic model parameters (e.g., temperature of 0.7), and not guiding their prompts to the AI system at all, besides requiring at least one interaction per question.

> 8 Broader Impacts and Ethical Considerations

Our work is driven by broader impacts: we seek to make AI evaluation more realistic and humancentered, by investigating how evaluation conclusions change when we incorporate human interactions. With our carefully designed user study, we show that evaluation conclusions change significantly from AI-alone to user-AI settings (for the same set of questions), and these results hold over different subject areas, AI models, AI-alone methods, and user-AI conditions. We hope that our work motivates AI researchers and practitioners to think more carefully about human-AI interactions when they evaluate AI systems, instead of only using AI-alone benchmarks. 692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

The direction of evaluating human-AI interactions also raises some ethical considerations. First, we should seek to recruit diverse human participants, since an AI system that works well for one individual or group may not work well for another (e.g., depending on ability, language, preferences, etc.). Second, user studies should be run ethically: participants should be paid fairly, they should provide informed consent about how their data will be used, and their data should be anonymized and personal information removed (e.g., if they tell the AI system their name). Third, the possibility of simulating humans in human-AI interactions is exciting and could make interactive evaluation feasible at scale, but LLM-based simulations of humans also have risks that need to be addressed, such as their possibilities for stereotyping, bias, and flattening populations (Cheng et al., 2023b,a; Wang et al., 2024). Researchers hoping to build and deploy user simulators should extensively probe for such biases, especially if user demographics are provided in simulator prompts.

References

- Lisa P. Argyle, Ethan C. Busby, Nancy Fulda, Joshua R. Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31:337–351.
- Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *CHI '21: Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, 81, pages 1–16.
- Alexander Bick, Adam Blandin, and David Deming. 2024. The rapid adoption of generative ai. *Federal Reserve Bank of St. Louis*.
- James Bisbee, Joshua D. Clinton, Cassy Dorff, Brenton Kenkel, and Jennifer M. Larson. 2024. Synthetic

743

744

- 797

replacements for human survey data? the perils of large language models. Political Analysis, 32:401-416.

- Serina Chang, Alicja Chaszczewicz, Emma Wang, Maya Josifovska, Emma Pierson, and Jure Leskovec. 2024. Llms generate structurally realistic social networks but overestimate political homophily. arXiv.
- Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023a. Marked personas: Using natural language prompts to measure stereotypes in language models. In ACL.
- Myra Cheng, Tiziano Piccardi, and Diyi Yang. 2023b. Compost: Characterizing and evaluating caricature in llm simulations. In *EMNLP*.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating llms by human preference. arXiv.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. arXiv.
- Katherine M. Collins, Albert Q. Jiang, Simon Frieder, Lionel Wong, Miri Zilka, Umang Bhatt, Thomas Lukasiewicz, Yuhuai Wu, Joshua B. Tenenbaum, William Hart, Timothy Gowers, Wenda Li, Adrian Weller, and Mateja Jamnik. 2024. Evaluating language models for mathematics through interactions. Proceedings of the National Academy of Sciences (PNAS).
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacafarm: A simulation framework for methods that learn from human feedback. In NeurIPS 2023.
- Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, Claire Barale, Robert McHardy, Joshua Harris, Jean Kaddour, Emile van Krieken, and Pasquale Minervini. 2024. Are we done with mmlu? *arXiv*.
- Gemini Team Google. 2023. Gemini: A family of highly capable multimodal models. arXiv.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In ICLR.
- Luke Hewitt, Ashwini Ashokkumar, Isaias Ghezae1, and Robb Willer. 2024. Predicting results of social science experiments using large language models. Working Paper.

John J. Horton. 2023. Large language models as simulated economic agents: What can we learn from homo silicus? arXiv.

798

799

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

- EunJeong Hwang, Bodhisattwa Majumder, and Niket Tandon. 2023. Aligning language models to user opinions. In EMNLP.
- Irina Jurenka, Markus Kunesch amd Kevin R. McKee, Daniel Gillick, Shaojian Zhu, Sara Wiltberger, Shubham Milind Phal, Katherine Hermann, Daniel Kasenberg, Avishkar Bhoopchand, et al. 2024. Towards responsible development of generative ai for education: An evaluation-driven approach. arXiv.
- Chuyi Kong, Yaxin Fan, Xiang Wan, Feng Jiang, and Benyou Wang. 2024. Platolm: Teaching llms in multi-round dialogue via a user simulator. In ACL.
- Mina Lee, Megha Srivastava, Amelia Hardy, John Thickstun, Esin Durmus, Ashwin Paranjape, Ines Gerard-Ursin, Xiang Lisa Li, Faisal Ladhak, Frieda Rong, Rose E. Wang, Minae Kwon, Joon Sung Park, Hancheng Cao, Tony Lee, Rishi Bommasani, Michael Bernstein, and Percy Liang. 2023. Evaluating human-language model interaction. Transactions on Machine Learning Research.
- Ruosen Li, Ruochen Li, Barry Wang, and Xinya Du. 2024a. IQA-EVAL: Automatic evaluation of humanmodel interactive question answering. In NeurIPS.
- Shuyue Stella Li, Vidhisha Balachandran, Shangbin Feng, Jonathan S. Ilgen, Emma Pierson, Pang Wei Koh, and Yulia Tsvetkov. 2024b. Mediq: Questionasking llms and a benchmark for reliable interactive clinical reasoning. In NeurIPS.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. 2024c. From live data to high-quality benchmarks: The arena-hard pipeline. LMSYS.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2023. Holistic evaluation of language models. Transactions of Machine Learning Research.
- Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Faeze Brahman, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. 2024. Wildbench: Benchmarking llms with challenging tasks from real users in the wild. In arXiv.
- Llama Team, AI@Meta. 2024. The llama 3 herd of models. arXiv.
- Evan Miller. 2024. Adding error bars to evals: A statistical approach to language model evaluations. arXiv preprint arXiv:2411.00640.

OpenAI. 2023. Gpt-4 technical report. arXiv.

Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In UIST.

Ruiyang Ren, Peng Qiu, Yingqi Qu, Jing Liu, Xin Zhao, Hua Wu, Ji-Rong Wen, and Haifeng Wang. 2024.
Bases: Large-scale web search user simulation with large language model based agents. In *EMNLP*.

853

857

865

870 871

872

873

875

876

878

897

901

- Hua Shen and Tongshuang Wu. 2023. Parachute: Evaluating interactive human-Im co-writing systems. In *CHI'23 In2Writing Workshop*.
- Michelle Vaccaro, Abdullah Almaatouq, and Thomas Malone. 2024. When combinations of humans and ai are useful: A systematic review and meta-analysis. *Nature Human Behaviour*, 8:2293–2303.
- Angelina Wang, Jamie Morgenstern, and John P. Dickerson. 2024. Large language models that replace human participants can harmfully misportray and flatten identity groups. *arXiv*.
 - Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *EMNLP*.
 - Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. Wildchat: 1m chatgpt interaction logs in the wild. In *ICLR*.
 - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging Ilm-as-a-judge with mt-bench and chatbot arena. In *NeurIPS*.

A Details on User Study

Task Details. We provide screenshots of all of the pages in our user study interface, including the Introduction Page (Figure A3), Phase 1 Tutorial (Figure A4), Confidence Page (Figure A5), User-Alone Page (Figure A6), Phase 2 Instructions (Figure A7), Phase 2 Tutorial (Figure A8), User-AI Page (Figure 2), and Feedback Page (Figure A9).

All Prolific workers who were located in the US, fluent in English, and had *not* participated in one of our pilots were eligible for our study. We used Prolific's standard sample, which distributed our study to available participants. Based on early pilots, we estimated that the task took around 25 minutes. We paid all participants \$5.00 upon completion of the entire task. We experimented with offering a small bonus per correct answer, which we discuss below. Our user study was approved by our institution's review board, which we will provide more details once anonymity is lifted.



Figure A1: Comparing results from Pilot 1 (without incentives) and Pilot 2 (with incentives).

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

A.1 Pilots and Incentivization

Pilot 1: no incentives. We ran one medium-sized pilot with 100 participants where we tested all datasets and models. At this point, we also included GPT-40-mini as a third model, in addition to GPT-40 and Llama-3.1-8b. In this pilot, we did not include incentives for correct answers. Results from this pilot did not show significant differences in accuracy between GPT-40 and GPT-40-mini, so we decided to drop GPT-40-mini from our full study, so that we could gather more answers per model.

Pilot 2: testing incentives. In our second pilot, we wanted to test the effect of including a small incentive for getting the correct answer, hypothesizing that it might improve the ecological validity of the study since users would try harder to answer the questions correctly. We included a small bonus of \$0.10 per correct answer, with a maximum bonus of \$1.30 for 13 questions, on top of the same base compensation of \$5.00 for completing the task.

While this bonus could help to improve ecological validity, there was a risk that the incentives result in users cheating on the study, such as by searching for the question on Google or ChatGPT. To mitigate this risk, first we repeatedly required users to acknowledge that they would not use external tools (Figures A3 and A7) and we said, "Compensation could be affected if we detect that you are using external tool." Second, we ran a second medium-sized pilot with incentives, with 60 participants on the three math datasets, and we compared the results between Pilots 1 and 2 to see if Pilot 2 had unrealistic increases in accuracy that could not be explained by slightly more user effort.

935

963

964

965

967

969

970

971

973

974

975

977

978

979

981

We visualize the mean accuracies per dataset and 936 model in Figure A1. We found that, as expected, in-937 centives tended to improve performance a little: out of 27 combinations of math datasets (3), models (3), 939 and answer types (i.e., user-alone, user-AI answerfirst, and user-AI direct-to-AI), the pilot with incen-941 tives had a higher mean accuracy 19 times. We also found that conversations were slightly longer with incentives. However, the overall improvement in accuracy was very small, only 3 percentage points, meaning we did not see unrealistic improvements 946 that would suggest use of external tools. We also 947 continued to see the gaps in user-AI performance 948 between the GPT models and Llama-3.1-8b, sug-949 gesting users were basing their answers on the AI Chatbot given to them. As further evidence of the 951 use of the AI Chatbot, and not external tools, we found that in the vast majority of cases (63 out 953 of 66 examples) where the user *changed* from an 954 incorrect user-alone answer to a correct user-AI answer, that new answer matched the answer given by the AI model in the user-AI conversation. Since 957 we found that incentives seemed to encourage users 958 959 to try slightly harder, and we did not see evidence of cheating, we decided to keep incentives for our 960 full study, but our pilot comparison shows that our 961 results were not overly sensitive to this decision.

A.2 ChatBench

In our full study, we recruited 650 participants and ran the study with incentives. Our dataset, Chat-Bench, compiles data over the full study and the two pilots. In Table A3, we provide additional data statistics, including how many answers we collected per model, dataset, condition, and answer type (user-alone or user-AI).

When releasing ChatBench, we will be careful to remove all personally identifying information (PII). We do not expect that there will be much PII, since the participants were clearly instructed to use the AI Chatbot to answer MMLU benchmark questions, not for their personal use, and we limited their number of messages per question and they had a time limit on the overall task. However, there is a risk, since participants could send any message they wanted to the AI Chatbot, so we will be careful to remove PII.

Dataset	Model	Comparison	Acc_1	SE_1	Acc_2	SE_2	<i>z</i> -value	p-value
Elementary Math	GPT-40	AI letter zero shot vs. UserAI direct to ai	0.73	0.04	0.92	0.02	-3.92	<0.01
Elementary Math	GPT-40	AI letter zero shot vs. UserAI answer first	0.73	0.04	0.90	0.02	-3.43	< 0.01
Elementary Math	GPT-40	AI letter few shot vs. UserAI direct to ai	0.74	0.04	0.92	0.02	-3.83	< 0.01
Elementary Math	GPT-40	AI letter few shot vs. UserAI answer first	0.74	0.04	0.90	0.02	-3.34	< 0.01
Elementary Math	GPT-40	AI free text vs. UserAI direct to ai	0.99	0.01	0.92	0.02	3.03	<0.01
Elementary Math	GPT-40	AI free text vs. UserAI answer first	0.99	0.01	0.90	0.02	4.04	<0.01
Elementary Math	GPT-40	User alone vs. UserAI direct to ai	0.78	0.03	0.92	0.02	-4.21	< 0.01
Elementary Math	GPT-40	User alone vs. UserAI answer first	0.78	0.03	0.90	0.02	-3.52	< 0.01
High School Math	GPT-40	AI letter zero shot vs. UserAI direct to ai	0.51	0.05	0.70	0.04	-3.20	< 0.01
High School Math	GPT-40	AI letter zero shot vs. UserAI answer first	0.51	0.05	0.73	0.03	-3.92	<0.01
High School Math	GPT-40	AI letter few shot vs. UserAI direct to ai	0.49	0.04	0.70	0.04	-3.57	<0.01
High School Math	GPT-40	AI letter few shot vs. UserAI answer first	0.49	0.04	0.73	0.03	-4.33	< 0.01
High School Math	GPT-40	AI free text vs. UserAI direct to ai	0.85	0.03	0.70	0.04	3.14	< 0.01
High School Math	GPT-40	AI free text vs. UserAI answer first	0.85	0.03	0.73	0.03	2.73	< 0.01
High School Math	GPT-40	User alone vs. UserAI direct to ai	0.41	0.03	0.70	0.04	-5.88	< 0.01
High School Math	GPT-40	User alone vs. UserAI answer first	0.41	0.03	0.73	0.03	-7.03	< 0.01
College Math	GPT-40	AI letter zero shot vs. UserAI direct to ai	0.45	0.07	0.73	0.05	-0.61	0.54
College Math	GPT-40	AI letter zero shot vs. UserAI answer first	0.45	0.07	0.52	0.07	-0.72	0.47
College Math	GPT-40	AI letter few shot vs. UserAI direct to ai	0.44	0.07	0.52	0.08	-0.72	0.47
College Math	GPT-40	AI letter few shot vs. UserAI answer first	0.44	0.07	0.52	0.07	-0.85	0.40
College Math	GPT-40	AI free text vs_UserAI direct to ai	0.73	0.06	0.52	0.08	2 23	0.03
College Math	GPT-40	AI free text vs. UserAI answer first	0.73	0.06	0.52	0.00	2.23	0.02
College Math	GPT-40	User alone vs. UserAI direct to ai	0.28	0.04	0.52	0.08	-2.67	< 0.01
College Math	GPT-40	User alone vs. UserAI answer first	0.28	0.04	0.52	0.07	-3.10	< 0.01
Conceptual Physics	GPT-40	AI letter zero shot vs. UserAI direct to ai	0.91	0.03	0.84	0.03	1.74	0.08
	CDT 4-		0.01	0.02	0.04	0.02	1.70	0.00
Conceptual Physics	GPT 40	AI letter few shot vs. UserAI direct to ai	0.91	0.03	0.84	0.05	1.70	0.09
Conceptual Physics	CPT 40	AI letter few shot vs. UserAI unect to al	0.90	0.02	0.04	0.03	3.22	< 0.01
Conceptual Physics	GPT 40	AI free text vs. UserAI direct to ai	0.90	0.02	0.84	0.03	3.22	
Conceptual Physics	GPT-40	AI free text vs. User AI answer first	0.97	0.02	0.84	0.03	3.62	
Conceptual I hysics	01 1-40	AT free text vs. OserAr answer first	0.97	0.02	0.04	0.05	5.05	<0.01
Conceptual Physics	GPT-40	User alone vs. UserAI direct to ai	0.55	0.03	0.84	0.03	-6.48	<0.01
Conceptual Physics	GPT-40	User alone vs. UserAI answer first	0.55	0.03	0.84	0.03	-6.69	< 0.01
Moral Scenarios	GPT-40	AI letter zero shot vs. UserAI direct to ai	0.71	0.05	0.79	0.03	-1.47	0.14
Moral Scenarios	GPT-40	AI letter zero shot vs. UserAI answer first	0.71	0.05	0.78	0.04	-1.13	0.26
Moral Scenarios	GPT-40	AI letter few shot vs. UserAI direct to ai	0.80	0.04	0.79	0.03	0.27	0.79
Moral Scenarios	GPT-40	AI letter few shot vs. UserAI answer first	0.80	0.04	0.78	0.04	0.49	0.63
Moral Scenarios	GPT-40	AI free text vs. UserAI direct to ai	0.72	0.05	0.79	0.03	-1.26	0.21
Moral Scenarios	GPT-40	AI free text vs. UserAI answer first	0.72	0.05	0.78	0.04	-0.93	0.35
Moral Scenarios	GPT-40	User alone vs. UserAI direct to ai	0.73	0.03	0.79	0.03	-1.54	0.12
Moral Scenarios	GPT-40	User alone vs. UserAI answer first	0.73	0.03	0.78	0.04	-1.05	0.29

Table A1: Results per dataset for GPT-40, including AI-alone vs. user-AI comparisons and user-alone vs. user-AI comparisons.

Dataset	Model	Comparison	Acc_1	SE_1	Acc_2	SE_2	<i>z</i> -value	p-value
Elementary Math	Llama-3 1-8b	AI letter zero shot vs. UserAI direct to ai	0.45	0.04	0.86	0.03	-8 58	< 0.01
Elementary Math	Llama-3.1-8b	AI letter zero shot vs. UserAI answer first	0.45	0.04	0.90	0.02	-10.50	< 0.01
Elementary Math	Llama-3.1-8b	AI letter few shot vs. UserAI direct to ai	0.43	0.03	0.86	0.03	-9.39	< 0.01
Elementary Math	Llama-3.1-8b	AI letter few shot vs. UserAI answer first	0.43	0.03	0.90	0.02	-11.53	< 0.01
Elementary Math	Llama-3.1-8b	AI free text vs. UserAI direct to ai	0.88	0.03	0.86	0.03	0.56	0.58
Elementary Math	Llama-3.1-8b	Al free text vs. UserAl answer first	0.88	0.03	0.90	0.02	-0.65	0.51
Elementary Math	Llama-3.1-8b	User alone vs. UserAI direct to ai	0.81	0.03	0.86	0.03	-1.26	0.21
Elementary Math	Llama-3.1-8b	User alone vs. UserAI answer first	0.81	0.03	0.90	0.02	-2.70	< 0.01
High School Math	Llama-3.1-8b	AI letter zero shot vs. UserAI direct to ai	0.32	0.03	0.62	0.04	-6.14	< 0.01
High School Math	Llama-3.1-8b	AI letter zero shot vs. UserAI answer first	0.32	0.03	0.64	0.04	-6.89	<0.01
High School Math	Llama-3.1-8b	AI letter few shot vs. UserAI direct to ai	0.30	0.02	0.62	0.04	-7.09	< 0.01
High School Math	Llama-3.1-8b	AI letter few shot vs. UserAI answer first	0.30	0.02	0.64	0.04	-7.98	< 0.01
High School Math	Llama-3.1-8b	AI free text vs. UserAI direct to ai	0.64	0.04	0.62	0.04	0.24	0.81
High School Math	Llama-3.1-8b	AI free text vs. UserAI answer first	0.64	0.04	0.64	0.04	-0.16	0.87
High School Math	Llama-3.1-8b	User alone vs. UserAI direct to ai	0.45	0.03	0.62	0.04	-3.37	<0.01
High School Math	Llama 3.1.8h	User alone vs. User AL answer first	0.45	0.03	0.64	0.04	3 03	<0.01
College Math	Liama-3.1-8b	AI letter zero shot vs. UserAI direct to ai	0.45	0.03	0.04	0.04	-3.93	0.17
College Math	Llama 2.1.8h	Al letter zero shot vs. UserAl answer first	0.35	0.04	0.40	0.07	-1.57	0.17
College Math	Liama 2.1.9h	Al letter for abot on UserAl direct to a	0.55	0.04	0.46	0.07	-1.30	0.12
College Math	Liama-3.1-8b	Al letter few shot vs. UserAl direct to al	0.30	0.04	0.46	0.07	-1.97	0.05
College Math	Liama-3.1-8b	Al letter few shot vs. UserAl answer first	0.30	0.04	0.48	0.07	-2.18	0.03
College Math	Llama-3.1-8b	AI free text vs. UserAI direct to ai	0.41	0.05	0.46	0.07	-0.57	0.57
College Math	Llama-3.1-8b	AI free text vs. UserAI answer first	0.41	0.05	0.48	0.07	-0.74	0.46
College Math	Llama-3.1-8b	User alone vs. UserAI direct to ai	0.40	0.04	0.46	0.07	-0.75	0.46
College Math	Llama-3.1-8b	User alone vs. UserAI answer first	0.40	0.04	0.48	0.07	-0.93	0.35
Conceptual Physics	Llama-3.1-8b	AI letter zero shot vs. UserAI direct to ai	0.53	0.05	0.67	0.04	-2.25	0.02
	11 21.01		0.52	0.05	0.72	0.04	2.00	.0.01
Conceptual Physics	Llama-3.1-8b	Al letter zero shot vs. UserAl answer first	0.53	0.05	0.73	0.04	-3.22	< 0.01
Conceptual Physics	Llama-3.1-8b	Al letter few shot vs. UserAl direct to al	0.57	0.04	0.67	0.04	-1.64	0.10
Conceptual Physics	Llama-3.1-8b	Al letter few shot vs. UserAl answer first	0.57	0.04	0.73	0.04	-2.70	< 0.01
Conceptual Physics	Llama-3.1-8b	Al free text vs. UserAl direct to ai	0.62	0.04	0.67	0.04	-0.77	0.44
Conceptual Physics	Llama-3.1-8b	AI free text vs. UserAI answer first	0.62	0.04	0.73	0.04	-1.80	0.07
Conceptual Physics	Llama-3.1-8b	User alone vs. UserAI direct to ai	0.46	0.03	0.67	0.04	-3.91	<0.01
Conceptual Physics	Llama-3.1-8b	User alone vs. UserAI answer first	0.46	0.03	0.73	0.04	-4.97	< 0.01
Moral Scenarios	Llama-3.1-8b	AI letter zero shot vs. UserAI direct to ai	0.40	0.03	0.72	0.04	-6.01	< 0.01
Moral Scenarios	Llama-3.1-8b	AI letter zero shot vs. UserAI answer first	0.40	0.03	0.74	0.04	-7.42	< 0.01
Moral Scenarios	Llama-3.1-8b	AI letter few shot vs. UserAI direct to ai	0.31	0.03	0.72	0.04	-7.35	< 0.01
M 10	11 210		0.01	0.02	0.74	0.04	0.06	.0.01
Moral Scenarios	Liama-3.1-8b	Al letter few shot vs. UserAl answer first	0.31	0.03	0.74	0.04	-8.86	< 0.01
Moral Scenarios	Liama-3.1-8b	AI free text vs. UserAI direct to at	0.49	0.03	0.72	0.04	-4.07	< 0.01
Moral Scenarios	Llama-3.1-8b	AI free text vs. UserAI answer first	0.49	0.03	0.74	0.04	-5.15	< 0.01
Moral Scenarios	Llama-3.1-8b	User alone vs. UserAl direct to ai	0.79	0.03	0.72	0.04	1.34	0.18
Moral Scenarios	Llama-3.1-8b	User alone vs. UserAI answer first	0.79	0.03	0.74	0.04	1.00	0.32

Table A2: Results per dataset for Llama-3.1-8b, including AI-alone vs. user-AI comparisons and user-alone vs. user-AI comparisons.

Model	Dataset	Condition	Answer Type	# Answers
gpt-40	college mathematics	answer-first	user-AI	134
01	C C		user-alone	283
		direct-to-AI	user-AI	116
			user-alone	121
	conceptual physics	answer-first	user-AI	318
			user-alone	425
		direct-to-AI	user-AI	352
			user-alone	117
	elementary mathematics	answer-first	user-AI	542
			user-alone	698
		direct-to-AI	user-AI	463
			user-alone	122
	high school mathematics	answer-first	user-AI	540
			user-alone	689
		direct-to-AI	user-AI	465
			user-alone	123
	moral scenarios	answer-first	user-AI	242
			user-alone	332
		direct-to-AI	user-AI	398
			user-alone	135
llama-3.1-8b	college mathematics	answer-first	user-Al	118
			user-alone	249
		direct-to-Al	user-Al	115
			user-alone	123
	conceptual physics	answer-first	user-Al	317
			user-alone	429
		direct-to-Al	user-Al	333
	<u> </u>		user-alone	112
	elementary mathematics	answer-first	user-Al	481
			user-alone	615
		direct-to-Al	user-Al	462
		<u> </u>	user-alone	123
	high school mathematics	answer-first	user-Al	475
		AT	user-alone	605
		direct-to-AI	user-Al	464
			user-alone	125
	moral scenarios	answer-nrst	user-AI	349
		diment to AT	user-alone	4/1
		direct-to-Al	user-AI	231
			user-alone	81

Table A3:	Dataset	statistics	for	ChatBer	ich.

Interactive Question Answering

INTRODUCTION

Thank you for taking the time to consider volunteering in a research project. This form explains what would happen if you join this research project. Please read it carefully and take as much time as you need. Email the study team to ask about anything that is not clear. Participation in this study is voluntary and you may withdraw at any time.

TITLE OF RESEARCH PROJECT

Interactive Question Answering

PRINCIPAL INVESTIGATOR

PURPOSE

The purpose of this project is to see how people answer questions when they have access to an AI tool.

PROCEDURES

During this project, you will be asked to try to answer multiple-choice questions. In Phase 1, you will try to answer questions on your own. In Phase 2, you will try to answer questions with the help of an AI chatbot. We will record your answers in both phases. In total, your participation will take around 25 minutes.

PERSONAL INFORMATION

Personal information we collect. Aside from your platform specific ID (e.g., Mechanical Turk ID etc.), no personal information will be collected during this study. Your platform specific ID can only be linked to your name by the platform, not by researchers, and the platform will not have access to your responses to this task. Your ID number will not be shared outside of _______ and the confines of this study without your permission, and will be promptly deleted after compensation has been successfully provided (30 days or less). De-identified data may be used for future research or given to another investigator for future use without additional consent. Researchers may share the results of this study publicly, such as in journal articles or conference presentations, but your name will not be included.

HOW YOU CAN ACCESS AND CONTROL YOUR INFORMATION

Once your platform specific ID is disassociated from your responses, we would not be able to remove your data from the study. For additional information or concerns about how handles your personal information, please see

BENEFITS AND RISKS

Benefits: There are no direct benefits to you that might reasonably be expected as a result of being in this study. The research team expects to learn about how AI tools help users answer questions, and there may be a public benefit to these research results being shared with the greater scientific community.

Risk: The risks of participating are similar to what you might experience while performing everyday tasks.

PAYMENT FOR PARTICIPATION

You will receive the compensation that was provided in the study description (\$5.00) after completing the entire study. You will also receive a small bonus of \$0.10 for each question you answer correctly, with a maximum total bonus of \$1.30 for 13 questions.

If you are unable to submit the study due to technical difficulties on your end, there is a risk of loss of payment. To mitigate, participants can reach out to the research team for input on resolving any difficulties encountered.

CONTACT INFORMATION

Should you have any questions concerning this project, please contact us at a research subject, please contact a should you have any questions about your rights as

CONSENT

By clicking "I agree" below, you confirm that the study was explained to you, you had a chance to ask questions before beginning the study, and all your questions were answered satisfactorily. By clicking "I agree" below, you voluntarily consent to participate, and you do not give up any legal rights you have as a study participant. If you would like a copy of this consent form, please print or save now. On behalf of we thank you for your contribution.

If you agree to participate, please click the continue button below. If you don't, please close this study.

I agree to this consent form.



Figure A2: Consent page. Parts are redacted to remain anonymous.

Introduction

This study contains two phases. In each phase you will be given a set of questions to answer. In **Phase 1**, you will be given **4 questions** to answer. In **Phase 2**, you will be given an additional **9 questions** to answer with the help of an AI chatbot.

For each of the 4 questions in Phase 1:

1. Read the question carefully.

2. Report your confidence in being able to answer the question before attempting to answer it.

3. Answer the question to the best of your ability, without using any external tools.

Please do not consult any external tools (e.g., ChatGPT, Google, Bing) while completing this study.

We're just interested in your best efforts and your experience. You will be paid regardless of how well you do!

The study will be ruined if you use external tools to do this task. Compensation could be affected if we detect that you are using external tools.

Please check below to indicate that you understand this and are ready to continue.

I promise not to use external tools to do this study, since it would ruin the study.

CONTINUE

Figure A3: Introduction page. Explains the task to users and ensures that they do not consult external tools.

Phase 1 Tutorial

This is an example of the type of question you will see in this study. To answer, select one of the multiple choice buttons, then press continue.

PHASE 1 Question 1 of 5
Now, answer this question to the best of your ability, without using any external tools. Take as much time as you need.
What is 4x3+5?
A. 17
O B. 20
○ C. 15
O D. 19
CONTINUE

Please check below to indicate that you understand this and are ready to continue.

□ I understand the instructions above and am ready to continue.

Figure A4: Phase 1 Tutorial. Provides an example of a Phase 1 question before the user begins Phase 1.

PHASE 1 Question 1 of 4

Please read the following question carefully.

A bag has 4 red and 6 blue marbles. A marble is selected and not replaced, then a second is selected. What is the probability that both are the same color?

\bigcirc	Α.	$\frac{1}{8}$
\bigcirc	В.	8 15
\bigcirc	C.	$\frac{1}{15}$
\bigcirc	D.	$\frac{7}{15}$

Before you begin, how confident are you that you know how to approach this problem?

- O Not confident: I don't know how to approach this problem.
- Somewhat confident: I have a rough idea of how to approach this problem but am not sure about the details.

O Very confident: I know how to approach this problem and am clear on the details.



Figure A5: Confidence page. Included per-question in both phases before the user tries to answer each question.

PHASE 1 Question 1 of 4

Now, answer this question to the best of your ability, without using any external tools. Take as much time as you need.

A bag has 4 red and 6 blue marbles. A marble is selected and not replaced, then a second is selected. What is the probability that both are the same color?

O	A. $\frac{1}{8}$
0	B. ⁸ / ₁₅
0	C. $\frac{1}{15}$
0	D. <u>7</u> 15
С	ONTINUE

Figure A6: User-alone page. Users answer all questions on their own in Phase 1 and, if they are in the *answer-first* condition, answer each question in Phase 2 on their own first before answering with AI.

Phase 2 Instructions

Thank you for completing Phase 1 of the study! In Phase 2, you will be given an additional 9 questions to answer.

- For each of the 9 questions in Phase 2:
- 1. Read the question carefully.

2. Report your confidence in being able to answer the question before attempting to answer it.

3. Answer the question to the best of your ability, without using any external tools

4. Answer the question again, this time using the AI chatbot provided in this study.

In Step 4, we expect you to use the AI chatbot. If you submit answers without sending any messages to the AI chatbot, your answers will be flagged and this could affect compensation.

Also, please only use the AI chatbot provided in this study and don't consult external tools (e.g., ChatGPT, Google, Bing).

We're just interested in your best efforts and your experience. You will be paid regardless of how well you do!

The study will be ruined if you use external tools to do this task. Compensation could be affected if we detect that you are using external tools.

Please check below to indicate that you understand this and are ready to continue.

I promise to use the AI chatbot in this study and not to use external tools.

CONTINUE

Figure A7: Phase 2 Instructions. Explains to users what they can expect in Phase 2 and reminds them not to use external tools. This screenshot shows instructions for a user in the *answer-first* condition. Users in the *direct-to-AI* condition see similar instructions, but without Step 3.

Phase 2 Tutorial

In Phase 2, you will have access to an AI Chatbot to help you with answering each question. Use the AI Chatbot as you would normally use any chatbot. You may send a maximum of 20 messages per question. As before, answer the question by selecting one of the multiple choice buttons, then press continue.

PHASE 2 Question 1 of 5	AI Chatbot
Now, try to answer the same question with the help of the AI Chatbot, using it as you would normally use any chatbot. You may say anything you want to the AI Chatbot to help you with answering the question.	
What is 7+(8-3)x2?	
○ A.21	
O B. 17	
O C. 20	
0.22	
CONTINUE	
	You can send 20 more messages to AJ.

Please check below to indicate that you understand this and are ready to continue.

□ I understand the instructions above and am ready to continue.

Figure A8: Phase 2 Tutorial. Provides an example of a Phase 2 question before the user begins Phase 2.

A few last questions

How difficult did you find it to answer these questions on your own? *

How useful was the chatbot for answering these questions? *

Did you notice any mistakes in the chatbot's responses? *

Did you have any technical issues during the study? Is there anything you would like us to know? *

CONTINUE

Figure A9: Feedback page. Final page of the task, where users leave free-text feedback to various questions.

1066

1067

1069

1070

1071

1073

1074

1075

1076

1077

1078

1079

1081

1082

1083

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1097

B Details on Analyses and Experiments

We download the MMLU datasets⁴ (Hendrycks et al., 2021) and MMLU-Redux datasets⁵ (Gema et al., 2024) from Hugging Face. The datasets are protected by the MIT and CC-by-4.0 licenses, respectively, allowing our use of this data in our research.

B.1 AI-Alone experiments

Here we provide the exact prompts used for each of the AI-alone methods: few-shot letter-only (Listing 1), zero-shot letter-only (Listing 2), and the two prompts for free-text (Listing 3 and Listing 4). For all the methods, the system prompt was "You are a helpful AI assistant."

Listing 1: Prompt for few-shot letter-only, taken from HELM. In-context examples are the five examples in MMLU's "dev" set for this dataset.

Answer with only a single letter.
The following are multiple choice questions (with answers) about {dataset}.
<pre>{example_1} A. {example_1_option_A} B. {example_1_option_B} C. {example_1_option_C} D. {example_1_option_D} Answer: {example_1_answer}</pre>
<pre>{example_5} A. {example_5_option_A} B. {example_5_option_B} C. {example_5_option_C} D. {example_5_option_D} Answer: {example_5_answer}</pre>
<pre>{question} A. {option_A} B. {option_B} C. {option_C} D. {option_D} Answer:</pre>

Listing 2: Prompt for zero-shot letter-only, using the same language as few-shot but dropping the in-context examples.

Answer with	only a	single	letter.
{question} A. {option_ B. {option_ C. {option_ D. {option_ Answer:	A } B } C } D }		

Listing 3: First prompt for AI-alone free-text. This prompt to generate the model's free-text response is simply the question and answer options concatenated.

{questi	on}		
A. {opt	ion_A}		
B. {opt	ion_B}		
C. {opt	ion_C}		
D. {opt	ion_D}		

⁴https://huggingface.co/datasets/cais/mmlu
⁵https://huggingface.co/datasets/
edinburgh-dawg/mmlu-redux-2.0

Listing 4: Second prompt for AI-alone free-text. This second prompt instructs GPT-40 to extract an answer (if any) from the model's free-text response. In order to not bias the answer extraction, we do not include the correct answer in this prompt.

B.2 Statistical details

Mean accuracies. When measuring accuracies for all methods (user-alone, AI-alone, and user-AI), we first compute per-question accuracies as the fraction of correct answers over total answers n_q for each question, denoted \hat{p}_q . We also compute the standard error for each question-level accuracy estimate $SE_q = \sqrt{\hat{p}_q(1-\hat{p}_q)/n_q}$. We then compute dataset-level accuracies with an (unweighted) average across all Q question-level accuracies, and dataset-level standard errors using decomposition of total variance to account for both variability in sampling questions from the larger population of MMLU questions and variability in correctness of responses: $SE_{tot} = \sqrt{(E[SE_q] + Var(\hat{p}_q))/Q}$ (Miller, 2024).

In Tables A1 and A2, we report mean accuracies for all datasets, models, AI-alone methods, and user-AI conditions. We also compare accuracies between two methods, for AI-alone vs. user-AI and for user-alone vs. user-AI. We conduct z-tests for all statistical tests comparing accuracies between two methods where

$$z = (\hat{p}_1 - \hat{p}_2) / \sqrt{SE_1^2 + SE_2^2}.$$
 (1)

Upper-bound on correlation. Since there is noise in our estimate of user-AI accuracy per question, we want to check if the low correlations between user-AI and AI-alone accuracies can be explained by that noise. To test this, we simulate an upper bound on what the correlation would be if the user-AI accuracies were drawn from the same distribution as the AI-alone accuracies, which we assume are perfectly estimated because we test each

982

985

988

990

991

992

LLM 50 times on each question. We construct hy-1098 pothetical user-AI data, where for each question q, 1099 we draw x from $Binomial(n_{qu}, p_{qf})$, where n_{qu} is 1100 the number of user-AI answers we had in our study, 1101 p_{qf} is the free-text accuracy on this question, and 1102 x/n_{qu} is now the hypothetical user-AI accuracy on 1103 this question. Over 100 iterations, the hypothetical 1104 correlation with free-text ranges from r = 0.881105 to 0.92 for *direct-to-AI* and r = 0.90 to 0.94 for 1106 answer-first, with the real correlations of 0.45 and 1107 0.46 (respectively) well below this range. 1108

B.3 Additional analyses from user study

User-AI conversations. To characterize the conversations between users and LLMs, we set a system prompt of "You are an expert at extracting information from conversations and MUST return a JSON object." and used the prompt in Listing 5.

Listing 5: Prompt to GPT-40 for automatically characterizing user-AI conversations.

The following conversation occurred between a user called You" and an AI called "Bot", in which "You" tried to use "Bot" to answer a question. Extract the following	"
information about this conversation. 1. Is the first substantive user prompt a: guestion.	
statement, or other?	
2. Is the first user question: a near-"exact" rephrasi	ng
of the question, a component of the question, ar	ı
erroneous or misinterpreted aspect of the questic	on,
or other?	
3. Is the first AI answer: correct, incorrect, or not	
yet providing an answer?	
 Is the last Al answer: correct, incorrect, or not y 	et
providing an answer?	
5. Anywhere in the course of the conversation does AT	
correct the user?	
6. Anywhere in the course of the conversation, does th	е
user correct the AI by selecting a different answ	ver
than what the AI recommended?	
Anything unusual or interesting about this	
interaction that you noticed?	
Here is the question the user was given:	
{duestion}	
Here are the answer choices to the question:	
{choices}	
Here is the correct answer to the question:	
{correctAnswer}	
Here is the conversation between the user ("You") and	
the Al ("Bot"):	
{conversation}	
Here is what the user selected as their answer.	
{selectedAnswer}	
(seree comment)	

We report our main findings from this analysis in the main text, where we showed in Figure 5 that only 34% of user-AI interactions "mirror" AI benchmarks. We also use the structured data produced by this analysis to measure how often the user corrects the AI model's mistake, by computing the fraction of user-AI interactions where the last AI answer in the conversation is wrong but the user still answered correctly (Figure B1). We find that users are much likelier to correct Llama-3.1-8b than GPT-4o, which helps to explain how some



Figure B1: Fraction of user-AI interactions where the last AI answer in the conversation is wrong but the user still answered correctly, by subject and model.

of the gap in the model's AI-alone performance is closed in the user-AI setting.

1166

1167

User confidence. In Figure B2, we visualize the 1168 relationship between user-reported confidence per 1169 question and their user-alone accuracy. First, over 1170 our five datasets, we find that users are most confi-1171 dent about Moral Scenarios, followed by Elemen-1172 tary Math, Conceptual Physics, High School Math, 1173 and College Math. The user selects their confidence 1174 from three options, as shown Figure A5), "not con-1175 fident", "somewhat confident", and "very confi-1176 dent". We find that users are well-calibrated within 1177 dataset: as their confidence increases, so does 1178 the mean accuracy. Users are less well-calibrated 1179 across datasets: for example, users who are very 1180 confident on a Conceptual Physics question slightly 1181 underperform those who are only somewhat confi-1182 dent on an Elementary Mathematics question. 1183

1109

1110

1111

1112

1113

1114



Figure B2: Distribution of confidence answers from users and mean user-alone accuracies per confidence answer.

B.4 Simulator details

1184

1185

1186

1187

1200

Below we provide the exact prompts for the twostep simulator (Listings 6-8) and the IQA-EVAL simulator from Li et al. (2024a) (Listing 9).

Listing 6: Two-step user simulator, system prompt for both tasks.

You	are a are	human trying	user to a	inter nswer	acti the	ng with followi	an ng d	AI ques	system, stion:	and	you
{que A. { B. { C. { D. {	option option option option option	<pre>} n_A } n_B } n_C } n_D }</pre>									

Listing 7: Two-step user simulator, user prompt for Task 1 (user refers to the role in the OpenAI API, not a real user).

Generate	the	first	prom	pt	you	would	l say	to	the	system	to	get
st	arte	d with	ansv	veri	ing	your	quest	ion	. Re	member	to	
wri	te e	xactly	as a	a re	eal	user	would					

Listing 8: Two-step user simulator, user prompt for Task 2 (user refers to the role in the OpenAI API, not a real user).

Listing 9: IQA-EVAL simulator, only has system prompt, following the original implementation.

You are mimicking a human. You are trying to choose the correct answer to the given
question. Please ask an assistant sub-questions for help approaching
answers.
In each turn, please only ask one sub-question to interact
with an assistant. In the sub-questions, please include
and necessary information, such as the question and
answer, please output "So, the answer is: A. B. C. or [
" "
{question}
A. {option_A}
B. {option_B}
C. {option_C}
YOU: {simulator prompt 1}
SYSTEM: {AI system response 1}
YOU: {simluator prompt k}
SYSTEM: {AI system response k}

In our simulator experiments, we fine-tune GPT-40 using Azure OpenAI Service. We use the default hyperparameters, with a batch size of 11 and 2 epochs. The training data contains 8,538 training examples (we describe in Section 5 how each user-AI conversation with k user utterances becomes k + 1 training examples for fine-tuning).



Figure B3: Scatter plot comparing different AI-alone and user simulator methods' abilities to predict user-AI accuracy, where the AI system is GPT-40. Pearson correlations are included in the plot titles.



Figure B4: Scatter plot comparing different AI-alone and user simulator methods' abilities to predict user-AI accuracy, where the AI system is Llama-3.1-8b. Pearson correlations are included in the plot titles.