

LLM Roleplay: Simulating Human-Chatbot Interaction

Hovhannes Tamoyan, Hendrik Schuff, and Iryna Gurevych

Ubiquitous Knowledge Processing Lab (UKP Lab)
Department of Computer Science and Hessian Center for AI (hessian.AI)
Technical University of Darmstadt
www.ukp.tu-darmstadt.de

Abstract

The development of chatbots requires collecting a large number of human-chatbot dialogues to reflect the breadth of users' sociodemographic backgrounds and conversational goals. However, the resource requirements to conduct the respective user studies can be prohibitively high and often only allow for a narrow analysis of specific dialogue goals and participant demographics. In this paper, we propose LLM Roleplay: a goal-oriented, persona-based method to automatically generate diverse multi-turn dialogues simulating human-chatbot interaction. LLM Roleplay can be applied to generate dialogues with any type of chatbot and uses large language models (LLMs) to play the role of textually described personas. To validate our method, we collect natural human-chatbot dialogues from different sociodemographic groups and conduct a user study to compare these with our generated dialogues. We evaluate the capabilities of state-of-the-art LLMs in maintaining a conversation during their embodiment of a specific persona and find that our method can simulate human-chatbot dialogues with a high indistinguishability rate.¹

1 Introduction

Collecting human-chatbot dialogues requires recruiting and managing a large number of human annotators, which can pose prohibitive obstacles to researchers who aim to develop conversational AI agents (i.e., chatbots). To circumvent the latter and the limitations of publicly available data, numerous methods employing chatbots to generate dialogues have been introduced lately (Xu et al., 2023b; Kim et al., 2023; Zhu et al., 2023; Ding et al., 2023; Svikhnushina and Pu, 2023; Zhao et al., 2024b). These methods can generate dialogues much faster and more cost-effectively while still approximating the quality and variety of human annotators. (Zhang et al., 2024).

¹<https://github.com/UKPLab/llm-roleplay>

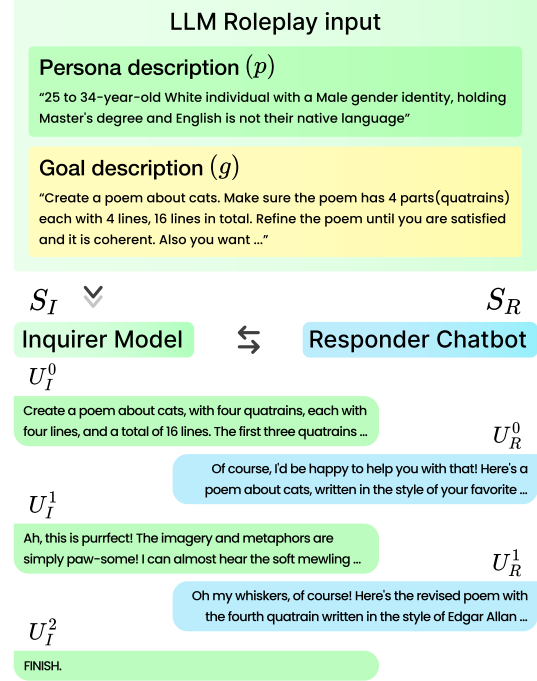


Figure 1: Schematic illustration of our method: A textual description of a persona and a goal (top) is used to instruct the inquirer (S_I) model to embody the given persona (left) and engage in a dialogue with the responder (S_R) chatbot (right). We show that dialogues simulated by the inquirer LLM and the responder chatbot can effectively simulate human-chatbot interaction.

Recently, several studies have leveraged such synthetic data generated through distillation and self-improvement techniques to enhance the capabilities of large language models (LLMs) (Zhang et al., 2024). By employing instruction tuning and fine-tuning techniques, these efforts aim to ensure that LLMs generate useful and safe responses that adhere to provided instructions (Xu et al., 2023a; Taori et al., 2023; Mukherjee et al., 2023; Li et al., 2023b; Peng et al., 2023; Almazrouei et al., 2023).

However, current dialogue generation methods have two critical limitations. First, the set of possible user responses for a given chatbot utterance

is large while the existing datasets typically only cover a single or few annotated user responses. Previously generated datasets are constrained by the number of dialogues and turns, making it difficult to extend them to novel domains or conversational goals. Second, existing methods rarely account for subjective response behavior and implicitly assume an average user. This can create a gap in representation (Rottger et al., 2022) and pose the risk of evaluating and improving chatbots for a specific sociodemographic group while neglecting others.

The quality of existing datasets is primarily assessed through model performance in a supervised fine-tuning (SFT) setup, which is influenced by various properties of the generated dataset. For example, Chen et al. (2023) argue that the initial prompt plays a crucial role in the quality of the generated dialogue. On the other hand, Zhao et al. (2024a) demonstrate that the length of the response is a key factor, and with significantly fewer samples, it outperforms methods that focus solely on quality. Meanwhile, Shen (2024) emphasizes the importance of mimicking human-style interactions in the dialogues.

To mitigate the shortcomings of existing dialogue generation methods and respond to the needs identified in prior work for SFT, we introduce LLM Roleplay: a goal-oriented, persona-based, plug-and-play method designed to generate diverse, multi-turn, long dialogues that simulate human-chatbot interactions using LLMs. Our method is the first to instruct an LLM (inquirer) to adopt a specific persona and prompt a chatbot (responder) to achieve a given conversational goal, thereby eliciting realistic human-AI interactions with a particular LLM. Figure 1 illustrates an exemplary application of our method.

In this work, we address the following two research questions: (i) To what extent can we simulate real human-chatbot dialogues using LLM-chatbot dialogues? (ii) How do various LLMs perform as inquirers within this setup? To answer those research questions, we conduct two user studies. First, we collect real human-chatbot dialogues by asking participants to reach various conversational goals and link the collected dialogues with participants’ sociodemographic backgrounds. Next, we use LLM Roleplay to simulate dialogues with the same set of personas and goals. Second, we conduct a human evaluation with another pool of participants to assess how well the generated dialogues mimic the collected ones. Specifically, we

present participants with two dialogues: one collected from the first study and one simulated using LLM Roleplay, both involving the same persona and the conversation goal.

We then ask the participants to identify which of the dialogues was simulated. We find that LLM Roleplay approximates natural human-chatbot dialogues with a high level of indistinguishability.

Overall, this paper contributes: (i) a novel method to simulate human-chatbot dialogues with an arbitrary set of personas and conversational goals; (ii) a human evaluation confirming our method’s potential to closely resemble real dialogues; (iii) a dataset of goal-oriented human-chatbot and model-chatbot dialogues using our hand-crafted multi-hop goals involving four state-of-the-art LLMs, (iv) an in-depth comparison of open-source and proprietary LLMs in maintaining conversations while embodying specific personas, (v) an open-source implementation of our plug-and-play method, readily applicable for simulating dialogues with any combination of models and chatbot systems across various conversational goals and personas.

2 Large Language Model Roleplay

In the following, we introduce the notation that we are following throughout this paper. Refer to Figure 1 for an annotated example.

Persona (\mathcal{P}) is the composition of sociodemographic features. We conceptualize persona as a written description of an individual from a specific sociodemographic group. We follow Kumar et al. (2021) to define the sociodemographic features and the options of the latter.

Goal (\mathcal{G}) is the textual representation of the conversational goal.

Subject (\mathcal{S}) is a dialogue participant. We denote the **inquirer** as $\mathcal{S}_{\mathcal{I}}$, the entity that asks questions, and the **responder** as $\mathcal{S}_{\mathcal{R}}$, the entity that answers the given questions. In our setup, the inquirer ($\mathcal{S}_{\mathcal{I}}$) is an LLM and the responder ($\mathcal{S}_{\mathcal{R}}$) is a conversational agent or chatbot (not necessarily an LLM). The human inquirer will be denoted as $\mathcal{S}_{\mathcal{I}}^h$.

Utterance (\mathcal{U}) is the output of a Subject (\mathcal{S}), i.e. either the output of the inquirer ($\mathcal{U}_{\mathcal{I}}$) or the responder ($\mathcal{U}_{\mathcal{R}}$). For example, $\mathcal{U}_{\mathcal{I}}^i$ will denote the inquirer’s i -th utterance. We refer to two consecutive utterances of two different subjects as a turn:

$$\mathcal{T}^i = [\mathcal{U}_I^i; \mathcal{U}_R^i].$$

Dialogue (\mathcal{D}) is a sequence of one or more turns. The dialogue of t turns is denoted as: $\mathcal{D}^t := \{\mathcal{T}^0, \mathcal{T}^1, \dots, \mathcal{T}^t\}$. We denote the maximum number of turns by max_t , i.e. $t \leq max_t$ which we discuss in more detail in Section 3.1.

At a high level, LLM Roleplay consists of three main steps: (i) Initial conditioning of the inquirer model with persona-specific text alongside the goal description. (ii) Subsequently, the inquirer’s output is provided to the responder model. (iii) Lastly, the output of the responder is returned to the inquirer by asking it to either output a follow-up question or a pre-defined token to terminate the dialogue.

We begin with assembling prompting templates for both of the subjects. We create a system prompt template (SYS_I) for the inquirer LLM. This template prompts the LLM to embody the given persona, provide a prompt to address the designated goal and output the specified termination token if it considers the goal accomplished. As for the responder, we use a default system prompt (SYS_R) to promote it to be a helpful and honest assistant. Additionally, we develop a response forwarder template (INTER_I) for the inquirer. This template prompts the inquirer to assess the conclusiveness of the answer of the responder, determining whether to output a subsequent question or the termination token. We provide the system and forwarder prompt templates in Table 5 in Appendix B. The algorithm for the LLM Roleplay is shown in Algorithm 1.

To obtain a dialogue for a given persona and a goal we create a prompt by passing the persona (\mathcal{P}) and the goal (\mathcal{G}) to the inquirer system prompt template and generate an output based on inquirer LLM distribution:

$$\mathcal{U}_I^0 \sim \mathcal{S}_I(\text{SYS_I}(\mathcal{P}, \mathcal{G})). \quad (1)$$

A deterministic prompt extraction function, `extract_prompt`, that looks for a string in double quotes in the response, is applied to the latter to extract the prompt from the response:

$$\mathcal{U}_I^0 := \text{extract_prompt}(\mathcal{U}_I^0). \quad (2)$$

The responder receives the prompt of the inquirer, and generates an output:

$$\mathcal{U}_R^0 \sim \mathcal{S}_R(\mathcal{U}_I^0). \quad (3)$$

Given the output of the responder, we condition the inquirer using the output forwarder template:

$$\mathcal{U}_I^1 \sim \mathcal{S}_I([\mathcal{U}_I^0; \text{INTER_I}(\mathcal{U}_R^0)]). \quad (4)$$

If the output of the inquirer begins or ends with the termination token, the stopping condition, `stop` is met, and consequently, the algorithm terminates. Otherwise, the prompt is extracted using `extract_prompt`, and the process persists for the coming turns until it reaches the maximum number of turns: $t = max_t$. For the t -th turn, the utterances will be:

$$\mathcal{U}_I^t \sim \mathcal{S}_I(\mathcal{T}^0, \mathcal{T}^1, \dots, \mathcal{T}^{t-1}) \quad (5)$$

$$\mathcal{U}_R^t \sim \mathcal{S}_R(\mathcal{T}^0, \mathcal{T}^1, \dots, \mathcal{T}^{t-1}, \mathcal{U}_I^t) \quad (6)$$

Algorithm 1: LLM Roleplay

Input : \mathcal{P}, \mathcal{G} : persona and goal
Output : $\mathcal{D}(\mathcal{S}_I, \mathcal{S}_R)$

```

1  $\mathcal{D} \leftarrow \{\}$ ;
2  $\mathcal{U}_I^0 \leftarrow \mathcal{S}_I(\text{SYS\_I}(\mathcal{P}, \mathcal{G}))$ 
3  $\mathcal{U}_I^0 := \text{extract\_prompt}(\mathcal{U}_I^0)$ 
4 if  $\mathcal{U}_I^0 = \emptyset$  then break;
5  $\mathcal{U}_R^0 \leftarrow \mathcal{S}_R(\mathcal{U}_I^0)$ 
6 for  $t = 1 \rightarrow max\_t$  do
7    $\mathcal{U}_I^t \leftarrow \mathcal{S}_I([\mathcal{U}_I^{t-1}; \text{INTER\_I}(\mathcal{U}_R^{t-1})])$ ;
8   if stop( $\mathcal{U}_I^t$ ) then
9     break;
10  end
11   $\mathcal{U}_I^t := \text{extract\_prompt}(\mathcal{U}_I^t)$ 
12  if  $\mathcal{U}_I^t = \emptyset$  then
13    break;
14  end
15   $\mathcal{U}_R^t \leftarrow \mathcal{S}_R(\mathcal{D}, \mathcal{U}_I^t)$ ;
16   $\mathcal{D} \leftarrow \{\mathcal{D}, [\mathcal{U}_I^t; \mathcal{U}_R^t]\}$ ;
17 end
```

3 Experiments

In the first study, we collect dialogues between human inquirers and model responders ($\mathcal{D}(\mathcal{S}_I^h, \mathcal{S}_R)$) by engaging participants with various personas (\mathcal{P}) in interactions with a chat-tuned LLM intending to achieve a given goal (\mathcal{G}). Utilizing the same set of personas (\mathcal{P}) and goals (\mathcal{G}) we generate dialogues between model inquirers and the same responder ($\mathcal{D}(\mathcal{S}_I, \mathcal{S}_R)$) by employing the LLM Roleplay method (Section 3.1). We report the statistics of the generated dialogues such as the average number of turns per dialogue, the number of tokens per

| | Llama-2 | Mixtral | Vicuna | GPT4 |
|---------------------------------|------------------------|----------------------|----------------------|----------------------|
| Avg. # Turns per Dialogue | 2.62 (1.54) | 3.86 (2.19) | 7.12 (3.79) | 7.60 (3.08) |
| Avg. # Tokens per Prompt | 77.77 (46.20) | 50.82 (26.47) | 75.19 (92.43) | 68.13 (60.37) |
| Avg. # Tokens per Response | 347.50 (142.14) | 302.47 (151.30) | 228.29 (163.92) | 267.69 (147.26) |
| No-prompt | 6.82% (1.48%) | 0.97% (0.57%) | 7.90% (0.54%) | 0.17% (0.04%) |
| Multiple Prompts | 8.79% (0.95%) | 8.40% (1.52%) | 6.08% (0.41%) | 39.21% (2.38) |
| Incoherent Response | 3.12% (0.30%) | 0.13% (0.09%) | 0.79% (0.10%) | 0.03% (0.04%) |
| Number of Self-Replies | 5.50% (3.25%) | 5.99% (1.43) | 69.39% (5.77%) | 5.48% (0.09%) |
| Incoherent Response (Responder) | 0.56% (0.57%) | 1.16% (0.19%) | 8.01% (0.93%) | 7.50% (0.35%) |

Table 1: Analysis of persona-specific dialogue collection (top) and failure cases (bottom) conducted for Llama-2, Mixtral, Vicuna, and GPT4. The results are averaged over runs with three different seeds. We show that GPT4 is successful at holding dialogues with longer utterances and having relatively fewer failure cases. Meanwhile, Mixtral is better at providing short on-point prompts. The number of utterances in dialogues is preferred to be larger, while for other metrics, smaller values are better. The standard deviation is indicated in parentheses.

prompt (inquirer output), and the number of tokens per response (responder output). Please refer to Appendix B for details on the generation parameters.

Furthermore, we investigate the failure cases of the inquirer LLMs, in following the given instruction. We depict the most common failure cases, report their statistics, and describe their detection mechanisms (Section 3.3).

In the second study, we conduct a human evaluation wherein another set of participants compare the natural $\mathcal{D}(\mathcal{S}_I^h, \mathcal{S}_R)$ and the simulated $\mathcal{D}(\mathcal{S}_I, \mathcal{S}_R)$ dialogues to discern the simulated counterpart $\mathcal{D}(\mathcal{S}_I, \mathcal{S}_R)$ (Section 3.4). We report the total and per-model undetectability rates, the utterance number on which the dialogue was detected, and the distribution of the duration and confidence choices.

Moreover, we use generalized linear models to analyze the detection rate of the simulated dialogue, the detection utterances number, and the duration users spent making a choice to analyze how different inquirer LLMs behave.

For our experiments, we used a single NVIDIA A100 GPU with 80GB memory for Llama-2 and Vicuna. We utilized up to 92% of the memory.

3.1 Persona-Specific Dialogue Collection

For this study, participants were instructed to interact with Llama-2² (\mathcal{S}_R) to accomplish a designated goal. We additionally ask participants to provide sociodemographic information (\mathcal{P}) such as age group, gender, race, level of education, and whether they identify as native English speaker. We base the choice of sociodemographic features and the options of the features on previous work by Kumar

et al. (2021). We provide a detailed screenshot of the persona information form interface in Figure 9, and a screenshot of our chat interface in Figure 10 in the Appendix B. In order to cover different conversational goals, we design ten handcrafted multi-hop goals (\mathcal{G}) spanning three domains: "Math", "Coding", and "General Knowledge".

We conduct a study involving 20 participants each engaged in tackling 10 goals, resulting in the generation of 200 natural human-chatbot interaction dialogues ($\mathcal{D}(\mathcal{S}_I^h, \mathcal{S}_R)$). We provide the full sociodemographic distribution of the participants in Figure 3 to Figure 7 in Appendix B.

Subsequently, we generate dialogues utilizing our LLM Roleplay method with a set of three state-of-the-art LLMs and one proprietary conversational agent as inquirers (\mathcal{S}_I): llama-2-13B-Chat (Llama-2) (Touvron et al., 2023), Mixtral-8x7B-Instruct-v0.1 (Mixtral) (Jiang et al., 2024), vicuna-13b-v1.5-16k (Vicuna) (Peng et al., 2023), and GPT4 (OpenAI et al., 2024). See sample natural and generated dialogues using Mixtral inquirer and Llama-2 responder from Figure 14 to Figure 18 in Appendix E.

In total, the dialogue collection and generation results in the creation of 200 natural human-chatbot ($\mathcal{D}(\mathcal{S}_I^h, \mathcal{S}_R)$) and 800 (with 4 inquirers) simulated ($\mathcal{D}(\mathcal{S}_I, \mathcal{S}_R)$) dialogue pairs.

We observe that, on average, GPT-4 (OpenAI et al., 2024) is successful at holding longer dialogues with 7.60 turns on average (Table 1). Mixtral (Jiang et al., 2024) on the other hand is better at generating shorter and on-point (based on manual analysis) prompts with 50.82 tokens per prompt. The generated dialogues have an average of 5.30 turns, each consisting of an average of 67.97 tokens per prompt, resulting in longer dialogues than those

²Due to limited deployment resources quantized Llama-2-13B-chat-GGUF is utilized (Touvron et al., 2023), employing the Q5_K_M quantization method with 5 bits (Frantar et al., 2023)

| SD Information | TTR | dist-1 | dist-2 |
|----------------------|----------------------|----------------------|----------------------|
| None | 0.281 (0.004) | 0.284 (0.003) | 0.625 (0.007) |
| Age | 0.572 (0.010) | 0.576 (0.010) | 0.892 (0.007) |
| Race | 0.582 (0.006) | 0.587 (0.005) | 0.880 (0.002) |
| Gender | 0.411 (0.003) | 0.415 (0.003) | 0.747 (0.002) |
| Education | 0.579 (0.015) | 0.578 (0.016) | 0.869 (0.005) |
| Is Native EN Speaker | 0.394 (0.003) | 0.394 (0.003) | 0.721 (0.005) |
| All | 0.606 (0.019) | 0.605 (0.019) | 0.872 (0.006) |

Table 2: Ablation study on the impact of sociodemographic (SD) information on the lexical diversity of the simulated user’s language measured via mean Type-Token Ratio (TTR) and Distinct-N (dist-n) along with corresponding variance reported in parentheses. We observe enhanced lexical diversity with the addition of individual sociodemographic features.

from the previous work (Table 4).

3.2 Impact of Persona-Specific Information

We assess how adding sociodemographic information to the persona description that is provided to the model affects the lexical diversity of the simulated user utterances. We compare three scenarios: (a) using a baseline prompt with no sociodemographic information in the persona description, (b) adding single sociodemographic features, (c) and including all features as described in our method. Table 2 shows the respective Type-Token Ratio (TTR) (Zipf, 2013) and Distinct-N (dist-n) (Li et al., 2016) values quantifying lexical diversity for Mixtral inquirer. We observe that adding single sociodemographic features as well as combined sociodemographic features to the model prompt significantly increases lexical diversity.

3.3 Failure Cases

Since the outputs of the LLMs are free-form, applying deterministic functions to their output proves challenging, resulting in a theoretically infinite number of turns. Therefore, there is a need for a set limit on the number of turns: *max_t*. Nevertheless, we try to capture the algorithm failure cases to have an automatic assessment of the inquirer LLM failure cases. We expect the inquirer LLM to provide the intended prompt enclosed in double quotes as we explicitly request this in our instructions. See a sample expected output and the failure cases examples in Table 8 in Appendix C. All of the LLMs in our experiments face the following issues.

Prompt not in double-quotes. The model fails to provide a prompt within double quotes, thus causing the dialogue to terminate.

Incoherent output. The responder produces a repetitive token sequence. We identify such outputs and preemptively end the dialogue. We utilize the incoherent function to analyze text for incoherent strings (see Algorithm 2 in Appendix C).

Incoherent output of responder. The inquirer model fails to detect the case when the responder outputs incoherent text, leading to unsuccessful dialogues. We spot such outputs and stop the dialog, using the same incoherent function.

Inquirer self-reply. The inquirer model fails to maintain its intended role and answers its own question, resulting in a fully generated dialogue in a single utterance. We detect this deterministically by examining the presence of any special tokens of the responder model within the output. For instance, in the case of Llama-2, this token appears as "[INST]", while for Vicuna, it manifests as "### Human:".

Multiple prompts. The inquirer outputs multiple strings enclosed in double quotes. As sometimes this overlaps with the previous case (inquirer self-reply), we select the first as the prompt.

Dialogue-stopping criterion failure. In the intermediate prompts, we ask the chatbot to output a pre-defined token when it "thinks" the goal assigned to it is achieved. Nonetheless, it follows a limited set of tokens; for instance, "FINISH" was utilized in our experiments.

We attribute these failures to two common issues in LLMs: limited context length and a restricted set of fine-tuned instructions (Kaddour et al., 2023).

Toxic content detection. To proactively address the generation of potentially harmful content, we employ the Llama-2 Guard model (Team, 2024) to filter out toxic dialogues.

In most cases, GPT4 outperforms other models with lower failure rates: 0.17% for responses without prompts, 0.03% for incoherent responses, and 5.48% for self-replies (Table 1). However, it struggles with providing a single prompt, failing 39.21% of the time. In contrast, Vicuna excels in generating single prompts, achieving a failure rate of 6.08%. Llama-2 receives fewer incoherent outputs, with 0.56%. Mixtral’s performance is intermediate, showing more balanced results across different metrics. Additional plots and detailed data are provided in Appendix B. In these experiments, we detected no unsafe content.

| Subset | | Llama-2 | Mixtral | Vicuna | GPT4 |
|------------|----------------------------------|------------------------|----------------------|----------------|--------------------|
| total | Undetectability Rate | 33.5% | 44.0% | 22.5% | 35.0% |
| | Confidence: "very confident" | 33 (16.50%) | 32 (16.00%) | 75 (37.50%) | 43 (21.50%) |
| | Confidence: "confident" | 108 (54.00%) | 83 (41.50%) | 75 (37.50%) | 97 (48.50%) |
| | Confidence: "somewhat confident" | 59 (42.50%) | 85 (25.00%) | 50 (20.00%) | 60 (30.00%) |
| detected | Duration | 86.10 (157.14) | 99.45 (98.81) | 94.63 (150.70) | 124.88 (227.36) |
| | Utterance Number | 1.72 (0.78) | 2.38 (1.35) | 2.50 (1.84) | 3.72 (2.35) |
| | Confidence: "very confident" | 24 (12.00%) | 19 (9.50%) | 65 (32.50%) | 32 (16.00%) |
| | Confidence: "confident" | 80 (40.00%) | 51 (25.50%) | 59 (29.50%) | 66 (33.00%) |
| undetected | Confidence: "somewhat confident" | 29 (14.50%) | 42 (21.00%) | 31 (15.50%) | 32 (16.00%) |
| | Duration | 120.44 (139.93) | 114.93 (157.46) | 62.24 (62.37) | 94.52 (124.03) |
| | Confidence: "very confident" | 9 (4.50%) | 13 (6.50%) | 10 (5.00%) | 11 (5.50%) |
| | Confidence: "confident" | 28 (14.00%) | 32 (16.00%) | 16 (8.00%) | 31 (15.50%) |
| | Confidence: "somewhat confident" | 30 (15.00%) | 43 (21.50%) | 19 (9.50%) | 28 (14.00%) |

Table 3: Analysis of human-evaluation results for detected, undetected, and total dialogues for Llama-2, Mixtral, Vicuna, and GPT4, showing confidence statistics as occurrences (percentages in parentheses). We exhibit that Mixtral has the highest undetectability rate of 44% (**50%** reflecting absolute indistinguishability). Moreover, it has the lowest confidence choice statistics for not confidently detected and the highest statistics for confidently undetected dialogues. GPT4 has the highest utterance number: 3.72, showing that it is identified in later utterances. The duration (in seconds) and the utterance number standard deviations are in parentheses.

3.4 Human-Evaluation

To answer the question of how well LLM inquirer and responder dialogues $\mathcal{D}(\mathcal{S}_I, \mathcal{S}_R)$ approximate dialogues between human inquirer and responder $\mathcal{D}(\mathcal{S}_I^h, \mathcal{S}_R)$, we conduct a human evaluation study. In each round, participants are shown two dialogues side by side, both featuring the same persona (\mathcal{P}) and solving the same goal (\mathcal{G}).

We allocated different sets of dialogue pairs for each participant, ensuring that no participant encounters multiple dialogues from the same user (from Section 3.1) solving the same goal. A new group of 20 participants was selected, with each participant tasked with reviewing 40 dialogue pairs. The participants selected for this study represent a wide range of occupational backgrounds. Some have no prior experience with chatbots, while others are industry professionals.

Participants are required to answer three questions for each dialogue pair: (i) Identify which dialogues they perceive as artificial (simulated): Choices include "1st (left)", "2nd (right)", or "Not sure" the latter is considered to be a tie. (ii) Express the level of confidence in their selection: Options are "Somewhat Confident," "Confident," and "Very Confident". (iii) Identify the specific utterance number that signifies the artificiality within the dialogue: Options depend on the number of utterances of the dialogue pairs. We provide a detailed view of the instructions provided to users in Figure 11 and the interfaces of the applications used for the study Figure 12 in Appendix D.

In conducting the human evaluation, we also track the amount of time participants take to respond to questions. This measure allows us to estimate a proxy measure of the complexity of the dialogue pairs. We assume that the longer it takes for a participant to make a decision, the more challenging the pair is, indicating that the simulated dialogue is difficult to discern.

The study shows that among the 800 samples, the simulated dialogues remained undetected in 33.75% (**50%** reflecting absolute indistinguishability) of the dialogues. Per model, statistics show that Mixtral has the highest undetectability rate of 44.0%, after which is GPT4 with 35.0% followed by Llama-2 and Vicuna with 33.5% and 22.5% respectively (Figure 2). See the per-confidence choice distribution in Figure 13 in Appendix D.

The highest utterance number on the detected set of dialogues has GPT4, meaning that it took more utterances for participants to recognize the simulated nature of the inquirer’s responses (Table 3). For the detected subset of dialogues, Mixtral has the highest percentages for all confidence choices: 9.50%, 25.50%, and 21.00% for "very confident", "confident" and "somewhat confident". Also, it is the best for the undetected pair of dialogues with 6.50%, 16.00%, and 21.50% respectively. Llama-2 excels in duration with 86.10 and 120.44 seconds for the detected and undetected dialogue pairs. However, it should be noted that the variation in duration is high, indicating that the participants did not complete the study at a consistent pace.

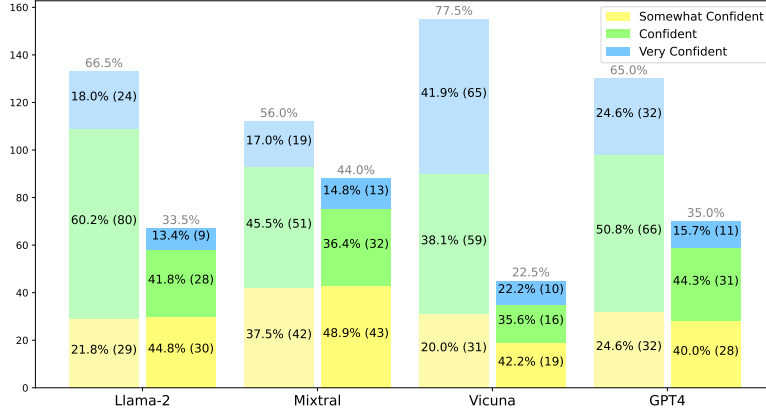


Figure 2: The distribution of detectability (left) and undetectability rates (right) per model for Llama-2, Mixtral, Vicuna, and GPT4. Each bar is stacked with confidence levels of: "Somewhat confident", "Confident" and "Very Confident". We shot that Mixtral has a relatively high undetectability rate of 44%, followed by GPT4 at 35%, Llama-2 at 33.5%, and Vicuna at 22.5%. The total (un)detectability rates for each model are mentioned in gray.

To further investigate how various instruction-tuned LLMs behave as inquirers within this setup we analyze the simulated dialogue detection probability, the detection utterance number, and the duration participants spent, using generalized linear mixed models (GLMMs), with the choice of LLM as the independent variable. We additionally include random effects to account for potential confounding effects of individual participants' detection abilities and users from Section 3.1.

We find significant effects of the choice of the model on the detection probability and the utterance position at which the participants formed their decision. We summarize the results of our statistical analyses in Table 10 in Appendix D using CLD codings (Piepho, 2004) and discuss our key findings in the following.

Effects on Detection Probability. We fit a binomial model (logit link) to predict the detection probability depending on the model. Concretely, we estimate a GLMM specified by: $detection_rate \sim model + (1|participant) + (1|generator_user)$. We observe a significant effect of the choice of the model on the detection probability ($\chi^2(3)=21.49$, $p < 0.001$). A post hoc Wald comparison of the contrasts for model types revealed significant differences between Vicuna and all other models (Vicuna being most likely to be detected), and Mixtral and Llama-2 (Mixtral being significantly less likely to be detected). We do not find a significant difference between GPT4 and Mixtral.

Effects on Utterance Positions. For the position of decision-forming utterances, we fit a respective

GLMM and find a significant effect of the choice of the model ($F=31.73$, $p < 0.001$). A post hoc Wald comparison of the contrasts for model types revealed significantly higher utterance positions for GPT4 than for the other models and that Llama-2 received significantly lower utterance position ratings. We do not observe a significant difference between Mixtral and Vicuna.

Our findings suggest that Mixtral, an open-source LLM, performs better than GPT4 in embodying a specified persona and simulating human-chatbot interactions in terms of detection probabilities. However, in detecting simulated dialogues based on utterance number, GPT4 outperforms, likely due to generating more utterances.

4 Related Work

Conversational Datasets. Approximation of the true distribution of human-chatbot dialogues is challenging and a significant number of diverse human annotators are needed. Incorporating participants from different sociodemographic profiles addressing a given conversational goal can substantially enhance the richness of the dataset. However, the associated time and resource requirements render the respectively needed large-scale studies infeasible for many researchers. Alternatively, a less resource-intensive and faster method for collecting human-chatbot conversations is to generate them using LLMs (Zhang et al., 2024). This method has been utilized in various setups for dialogue generation, such as human-human (Adiwardana et al., 2020; Kim et al., 2023; Chen et al., 2023; Li et al., 2023a), teacher-student (Macina et al., 2023), and

| Type | Dataset Name | # Dialogues | Avg. # Turns/Dialogue | Avg. # Tokens/Prompt | Avg. # Tokens/Response | Topics | Personalized |
|-----------|------------------------------------|-------------|-----------------------|----------------------|------------------------|-----------------|--------------|
| natural | OpenAssistant (Köpf et al., 2023) | 3k | 2.12 | 28.28 | 171.34 | human-crafted | yes |
| | LMSYS-Chat-1M (Zheng et al., 2024) | 777k | 1.92 | 55.23 | 163.66 | human-crafted | yes |
| | WildChat (Zhao et al., 2024b) | 360k | 2.46 | 164.32 | 276.50 | open | no |
| synthetic | UltraChat (Ding et al., 2023) | 1.5M | 3.85 | 52.54 | 249.41 | model-generated | no |
| | LLM Roleplay (Ours) | any | 5.30 (2.11) | 67.97 (10.51) | 286.48 (151.15) | any | yes |

Table 4: Statistics of key human-crafted (natural) and model-generated (synthetic) datasets (English subsets) relevant to our study. Our method can generate unlimited dialogues across various topics, featuring persona-based prompts and longer utterances. The natural datasets include personalized information reflecting the inquirer’s subjectivity. Standard deviations are shown in parentheses. For a comprehensive dataset list, see Table 9 in Appendix F.

patient-physician (Wang et al., 2023).

Prior work proposed a large number of human-crafted and synthetically generated datasets, each trying to collect more dialogue pairs revolving around various topics (Table 4). OpenAssistant (Köpf et al., 2023) collects human-crafted conversational dialogue trees, with prompts and answers generated by different humans. LMSYS-Chat-1M (Zheng et al., 2024) contains the refined version of dialogue logs collected via an online chat interface and predominantly contains dialogues between users and Vicuna-13b (Peng et al., 2023). UltraChat (Ding et al., 2023) and GLAN (Li et al., 2024) synthetically generate dialogues using proprietary conversational AI systems (ChatGPT Turbo) evolving around topics generated by the same systems.

Personas in Large Language Models. Large language models (LLMs) have demonstrated distinct behaviors and personas (Andreas, 2022; Wolf et al., 2024). Andreas (2022) note that LLMs interpret behaviors from text prompts, influencing generated content. Additionally, Wolf et al. (2024) argue that LLMs function as mixture decompositions, where prompts shift component balances, thus triggering persona-specific responses. Furthermore, Beck et al. (2024) investigate the effects of sociodemographic prompts on model responses, finding that while beneficial in some settings, they produce varied effects across models.

Building on these findings, our work introduces a novel, goal-oriented, persona-centric method for generating diverse, multi-turn dialogues. This method simulates dialogues across various combinations of conversational goals, personas, and LLMs, enabling the creation of countless simulated dialogues without theoretical limits.

5 Discussion and Future Work

In this work, we propose the novel LLM Roleplay method: an automatic, model-agnostic approach

for eliciting multi-turn, goal-oriented, persona-based simulated human-chatbot dialogues. We develop and validate our method through two user studies involving 40 participants. Our findings show that up to 44% of these dialogues, where 50% represents perfect indistinguishability, are indistinguishable from real human-chatbot interactions and feature more turns than previous datasets.

Building on this work’s findings, several future research avenues warrant exploration. First, our method can generate user-specific conversational datasets for targeted alignment and domain-adaptation (e.g., using RLHF (Christiano et al., 2017)). Second, it can improve dialogue evaluation by providing ample realistic conversational data.

While this paper presents important findings on LLM Roleplay and provides the first evidence that our method can approximate human-chatbot dialogues, the sociodemographic representativeness of our method still needs assessment and advancement. We release our dialogue dataset, method code, and synthesized dialogues to support future research in this promising field.

6 Conclusion

We present our novel LLM Roleplay method for simulating human-chatbot interaction. In a series of two user studies, we collect real human-chatbot dialogues and demonstrate that LLM Roleplay can generate diverse multi-turn conversations that approximate natural human-chatbot dialogues with a high level of indistinguishability. Our findings highlight the potential of LLMs in simulating human-chatbot interactions to synthesize realistic dialogues that create new opportunities for real-time model evaluation and training data generation for model fine-tuning.

7 Acknowledgments

This work has been funded by the German Research Foundation (DFG) as part of the UKP-

SQuARE project (grant GU 798/29-1). This work has been funded by the LOEWE Distinguished Chair “Ubiquitous Knowledge Processing”, LOEWE initiative, Hesse, Germany (Grant Number: LOEWE/4a//519/05/00.002(0002)/81). We gratefully acknowledge the support of Microsoft with a grant for access to OpenAI GPT models via the Azure cloud (Accelerate Foundation Model Academic Research).

Limitations

In this section, we explore the inherent limitations of this research study. We further note that all our experiments have been approved by the local ethics reviewing board at Technical University Darmstadt.

Sociodemographic Representativeness. An important limitation of our study lies in the sociodemographic distribution of the participants involved in the dialogue collection process. The pool of our studies’ participants serves as an initial investigation into the promising direction of interaction simulation and cannot represent the full spectrum of sociodemographic backgrounds. To address this limitation, future research needs to broaden the scope of participants’ sociodemographic groups and assess the replicability of our findings within large user studies, encompassing a more comprehensive range of sociodemographic groups. Studies should be conducted with specific sociodemographic groups, ensuring that participants representing a broader set of combinations of the persona features described in the paper are covered (e.g. via crowdsourcing on platforms like Prolific). By doing so, a more nuanced understanding of natural dialogue dynamics across diverse populations can be achieved, and allow us to uncover weaknesses and respectively needed improvements building upon our initial method specification.

Controllability. In this work, we have endeavored to detect and prevent failure cases of the proposed method. Despite our efforts, it is important to acknowledge that achieving absolute coverage in detecting all potential failure cases remains elusive. For example, instances where multiple turns of appreciation occur bilaterally present a challenge that we have not fully addressed. Moreover, we strive to save all detected failure cases by re-generating with different parameters. However, this approach has not proven to be effective. Future research should concentrate on exploring these scenarios further to

detect and prevent failure cases more efficiently.

Additionally, although our LLM Roleplay shows promising results in our experimental settings, it is not immune to the common challenges associated with large language models. Problems like hallucinations and associated LLM behavior issues can still arise, even though we have not encountered any during our experiments.

Ethics Statement

As discussed in the previous section, the narrow sociodemographic spectrum of participants involved in our user studies demands follow-up work to study the generated dialogues’ sociodemographic validity. Further, the sociodemographic information provided in our method’s prompts could potentially trigger biased content generation within the underlying LLM. As a first countermeasure, our method comprises a guard model to prevent the generation of toxic content. We, however, note that biases can also manifest more subtly and want to emphasize that, while we did not observe any such cases, future work should carefully assess whether such effects are present in future LLMs.

While synthetic data generation always entails a risk of generating invalid or biased data, we argue that our work takes an important step toward more valid user data generation. All our user studies have been approved by the local ethics reviewing board.

References

- Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. 2020. [Towards a human-like open-domain chatbot](#). *Preprint*, arXiv:2001.09977.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, M  rouane Debbah,   tienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [The falcon series of open language models](#). *Preprint*, arXiv:2311.16867.
- Jacob Andreas. 2022. [Language models as agent models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5769–5779, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Tilman Beck, Hendrik Schuff, Anne Lauscher, and Iryna Gurevych. 2024. [Sensitivity, performance, robustness: Deconstructing the effect of sociodemographic prompting](#). In *Proceedings of the 18th Conference of*

- the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2589–2615, St. Julian’s, Malta. Association for Computational Linguistics.
- Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Seokhwan Kim, Andy Rosenbaum, Yang Liu, Zhou Yu, and Dilek Hakkani-Tur. 2023. [PLACES: Prompting language models for social conversation synthesis](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 844–868, Dubrovnik, Croatia. Association for Computational Linguistics.
- Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. [Deep reinforcement learning from human preferences](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4299–4307.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. [Enhancing chat language models by scaling high-quality instructional conversations](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3029–3051, Singapore. Association for Computational Linguistics.
- Elias Frantar, Saleh Ashkboos, Torsten Hoeftler, and Dan Alistarh. 2023. [Gptq: Accurate post-training quantization for generative pre-trained transformers](#). *Preprint*, arXiv:2210.17323.
- Karthik Gopalakrishnan, Behnam Hedayatnia, Qinglang Chen, Anna Gottardi, Sanjeev Kwatra, Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-Tür. 2019. [Topical-chat: Towards knowledge-grounded open-domain conversations](#). In *Interspeech 2019, 20th Annual Conference of the International Speech Communication Association, Graz, Austria, 15-19 September 2019*, pages 1891–1895. ISCA.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio Cesar Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Conti Kauffmann, Gustavo Henrique de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Behl, Xin Wang, Sebastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. 2024. [Textbooks are all you need](#).
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2024. [Mixtral of experts](#). *Preprint*, arXiv:2401.04088.
- Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. [Challenges and applications of large language models](#). *Preprint*, arXiv:2307.10169.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2023. [SODA: Million-scale dialogue distillation with social commonsense contextualization](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12930–12949, Singapore. Association for Computational Linguistics.
- Andreas K  pf, Yannic Kilcher, Dimitri von R  tte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Minh Nguyen, Oliver Stanley, R  chard Nagyfi, Shahul ES, Sameer Suri, David Alexandrovich Glushkov, Arnav Varma Danturi, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Julian Mattick. 2023. [Openassistant conversations - democratizing large language model alignment](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Deepak Kumar, Patrick Gage Kelley, Sunny Consolvo, Joshua Mason, Elie Bursztein, Zakir Durumeric, Kurt Thomas, and Michael Bailey. 2021. Designing toxic content classification for a diversity of perspectives. In *Proceedings of the Seventeenth USENIX Conference on Usable Privacy and Security*, SOUPS’21, USA. USENIX Association.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023a. [CAMEL: Communicative agents for “mind” exploration of large language model society](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Haoran Li, Qingxiu Dong, Zhengyang Tang, Chaojun Wang, Xingxing Zhang, Haoyang Huang, Shaohan Huang, Xiaolong Huang, Zeqiang Huang, Dongdong Zhang, Yuxian Gu, Xin Cheng, Xun Wang, Si-Qing Chen, Li Dong, Wei Lu, Zhifang Sui, Benyou Wang, Wai Lam, and Furu Wei. 2024. [Synthetic data \(almost\) from scratch: Generalized instruction tuning for language models](#). *ArXiv preprint*, abs/2402.13064.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. [A diversity-promoting objective function for neural conversation models](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. [DailyDialog: A manually labelled multi-turn dialogue dataset](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*,

- pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023b. [Textbooks are all you need ii: phi-1.5 technical report](#). *Preprint*, arXiv:2309.05463.
- Jakub Macina, Nico Daheim, Sankalan Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023. [MathDial: A dialogue tutoring dataset with rich pedagogical properties grounded in math reasoning problems](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5602–5621, Singapore. Association for Computational Linguistics.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. [Orca: Progressive learning from complex explanation traces of gpt-4](#). *Preprint*, arXiv:2306.02707.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambatista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. [Instruction tuning with gpt-4](#). *Preprint*, arXiv:2304.03277.
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. [Red teaming language models with language models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3419–3448, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Hans-Peter Piepho. 2004. [An algorithm for a letter-based representation of all-pairwise comparisons](#). *Journal of Computational and Graphical Statistics*, 13:456–466.

- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. [Towards empathetic open-domain conversation models: A new benchmark and dataset](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Paul Rottger, Bertie Vidgen, Dirk Hovy, and Janet Pierrehumbert. 2022. [Two contrasting data annotation paradigms for subjective NLP tasks](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 175–190, Seattle, United States. Association for Computational Linguistics.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. [Character-LLM: A trainable agent for role-playing](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13153–13187, Singapore. Association for Computational Linguistics.
- Ming Shen. 2024. [Rethinking data selection for supervised fine-tuning](#). *ArXiv preprint*, abs/2402.06094.
- Ekaterina Svikhnushina and Pearl Pu. 2023. [Approximating online human evaluation of social chatbots with prompting](#). In *Proceedings of the 24th Meeting of the Special Interest Group on Discourse and Dialogue*, pages 268–281, Prague, Czechia. Association for Computational Linguistics.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7.
- Llama Team. 2024. Meta llama guard 2. https://github.com/meta-llama/PurpleLlama/blob/main/Llama-Guard2/MODEL_CARD.md.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Junda Wang, Zonghai Yao, Zhichao Yang, Huixue Zhou, Rumeng Li, Xun Wang, Yucheng Xu, and Hong Yu. 2023. [Notechat: A dataset of synthetic doctor-patient conversations conditioned on clinical notes](#). *Preprint*, arXiv:2310.15959.
- Yotam Wolf, Noam Wies, Oshri Avnery, Yoav Levine, and Amnon Shashua. 2024. [Fundamental limitation of alignment in large language models](#).
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. [Wizardlm: Empowering large language models to follow complex instructions](#). *Preprint*, arXiv:2304.12244.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023b. [Baize: An open-source chat model with parameter-efficient tuning on self-chat data](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6268–6278, Singapore. Association for Computational Linguistics.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2024. [Instruction tuning for large language models: A survey](#). *Preprint*, arXiv:2308.10792.
- Hao Zhao, Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. 2024a. [Long is more for alignment: A simple but tough-to-beat baseline for instruction fine-tuning](#). *ArXiv preprint*, abs/2402.04833.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024b. [\(in\)the wild chat: 570k chatGPT interaction logs in the wild](#). In *The Twelfth International Conference on Learning Representations*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhaghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2024. [LMSYS-chat-1m: A large-scale real-world LLM conversation dataset](#). In *The Twelfth International Conference on Learning Representations*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhaghao Wu, Yonghao Zhuang, Zi Lin,

Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging LLM-as-a-judge with MT-bench and chatbot arena](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. 2023. Starling-7b: Improving llm helpfulness & harmlessness with rlaiif.

George Kingsley Zipf. 2013. *The psycho-biology of language: An introduction to dynamic philology*. Routledge.

A Hand-Crafted Conversational Goals

In our preliminary experiments, we used one-hop conversational goals, each containing a singular question. We noticed a pattern where the inquirer replicates the primary question with subtle modifications and then presents it as the prompt. To make the interaction more intriguing and longer, we hand-craft multi-hop conversational goals. Specifically, 10 goals from three categories:

Math

- "You want to know how fast you run different distances. You use a stopwatch to measure the time it takes you to complete a 50-meter, 100-meter, and 200-meter race. You want to know how can you calculate your speed for each race? Based on that, you also want to calculate how many calories you burned during each race."
- "You can run at a rate of speed four times faster than you can walk, but you can skip at a rate of speed that is half as fast as you can run. You want to know If you can skip at 3 miles per hour, and how many miles can you travel in six hours if you spend one-third of the time and two-thirds of the time running and walking, respectively. Also, you are curious about the other way around (one-third of the time walking and two-thirds for running)."
- "Every day, you feed each of your chickens three cups of mixed chicken feed, containing seeds, mealworms, and vegetables to help keep them healthy. You give the chickens their feed in three separate meals. In the morning, you give your flock of chickens 15 cups of feed. In the afternoon, you give your chickens another 25 cups of feed. You want to know how many cups of feed you need to give your chickens in the final meal of the day if the size of your flock is 20 chickens. Also, you want to know how much the chicken egg production rate depends on the feed you give, and if you provide enough feed to your chickens for high-rate egg production."

Coding

- "You want to make this function better. You want the chatbot to make it recursive to have memory optimal function, but make sure that it doesn't enter into an infinite loop. After

that, you want to plug a CLI (command line interface) into this function, so the user can insert a number and get the factorial of it as output: 'The factorial of the <NUMBER>, is <FACTORIAL>'. "" def factorialize(num): factorial = 1 for i in range(1, num): factorial *= i return factorial ""

- "You have a little project where you need to use JavaScript, a language you don't use every day. You have a subtask to write a function that counts how many vowels are in a given string. And you need this functionality in OOP. Also, you want the chatbot to develop the snippet it provided by getting the function input string via an API call. If the chatbot uses functions or operators you are not familiar with feel free to ask follow-up questions about it."
- "You want to draw a unicorn in Python using the 'turtle' module. (There should be multiple lines of short function calls). After that substitute the 10th line, which includes number argument(s), with the value 73(s)."

General Knowledge

- "You want to know what are the world's 10 oldest continuously inhabited cities. Pick the 3rd in that list find out who established the city, in which region it is located and what was the highest population."
- "You have written content that disagrees with the following statement: 'Technology is the cause of all societal problems' And you want the chatbot to generate a response that agrees with the statement, to make your claims stronger."
- "You plan a trip to France and would like to do a walking tour. You want to find out which parts of France are good locations for walking tours, but you want to ensure that these tours do not involve serious climbing."
- "You want to use the chatbot to create a poem about cats. Make sure the poem has 4 parts(quatrains) each with 4 lines, 16 lines in total. Refine the poem until you are satisfied and it is coherent. Also, you want to change the style of one of the quatrains to reflect the distinctive style of your favourite poet."

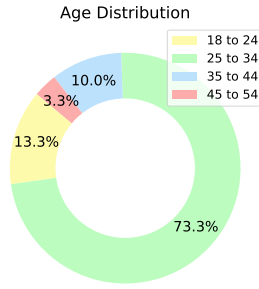


Figure 3: Age distribution of participants for persona-specific dialogue collection study

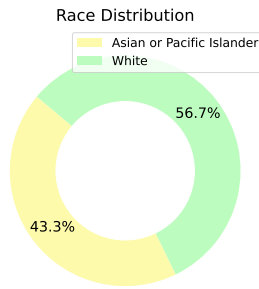


Figure 4: Race distribution of participants for persona-specific dialogue collection study

B Persona-specific Dialogues Collection

For persona-specific dialogue collection, we conducted a human study where participants were given the following instruction: "Below is your defined goal. What will you prompt the chatbot to accomplish your goal? Feel free to ask follow-up questions on the related topic of the question and clarify things in the response."

Participants for this study were selected from different sociodemographic groups, including individuals from four age groups "18 to 24", "25 to 34", "35 to 44", and "45 to 54", with "Asian or Pacific Islander" and "White" races, encompassing "female" and "male" genders, holding "Doctoral" and "Master's" degrees, and being either "native" or "non-native" English speakers. See the distributions of participants by features from Figure 3 to Figure 7.

Our initial experiments included falcon-40b-instruct (Almazrouei et al., 2023); however, we excluded it due to its difficulty in following instructions.

We use the default generation settings for all our models. By setting "do_sample=true" in the Hugging Face Transformers "generate()" method,

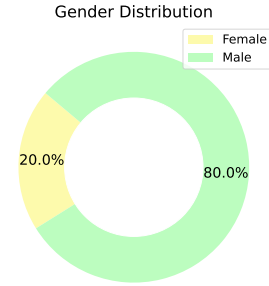


Figure 5: Gender distribution of participants for persona-specific dialogue collection study

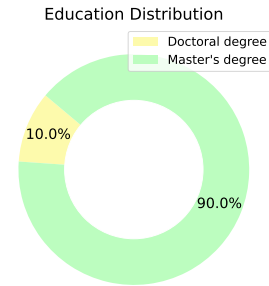


Figure 6: Education distribution of participants for persona-specific dialogue collection study

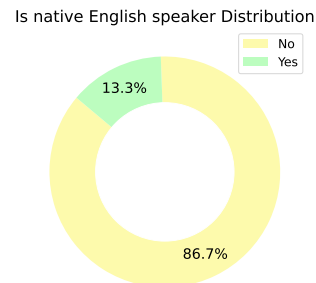


Figure 7: Is native English speaker distribution of participants for persona-specific dialogue collection study

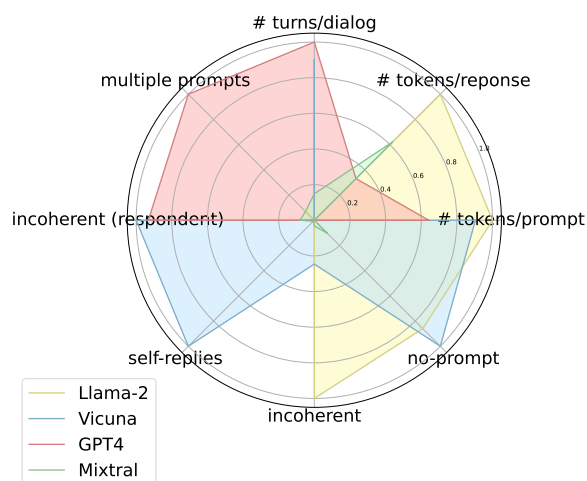


Figure 8: Normalized statistics of Llama-2, Mixtral, Vicuna and GPT4 for dialogue collection. The smaller the area of the model plot the better excluding the "# turns/dialog".

we enable multinomial sampling. For the inquirer model, we set 'max_new_tokens' to 1k, and for the responder model, we set it to 4k.

For Llama-2 and Vicuna inquirers we have used a single NVIDIA A100 GPU, however for Mixtral, we employed two NVIDIA A100 GPUs, with memory usage reaching a maximum of 92% and 66%, respectively. The experiments for Llama-2, Vicuna, Mixtral, and GPT4 inquirers took 2 hours and 45 minutes, 13 hours and 25 minutes, 17 hours and 38 minutes, and 9 hours and 6 minutes, respectively.

C Failure Cases

The incoherent text detection function considers n-grams up to a specified maximum, detects consecutive repetitions, and iterates through the text to examine increasing n-gram sizes. The function checks for repetitive patterns surpassing a specified threshold. If it finds such patterns, it returns True; otherwise, it returns False. See Algorithm 2.

In our experiments for Llama-2, the parameters `incoherent_max_n` and `incoherent_r` are set to 8 and 2 respectively. In the case of the Vicuna, these values are 5 and 2. For Mixtral and GPT4 they are 4 and 2.

D Human Evaluation

For the human evaluation we give the following instruction to the participants "Here, you'll find two dialogues: one is a conversation between a human and an AI, and the other is between AI and AI. Choose the dialogue

Figure 9: The screenshot of the person form from the application used for conducting the persona-specific dialogue collection containing the following fields: "Age Range", "Race", "Gender", "Education" and "Is English your native language"

Figure 10: The screenshot of the chat interface from the application used for conducting the persona-specific dialogue collection containing a simple chat interface with a "Next Goal" button to finish the current conversation and get the next conversational goal.

| Llama-2 | |
|---------------------|---|
| SYS_I | [INST] «SYS» You are <PERSONA>. You will start a conversation with an assistant. If you accomplish your ultimate goal during the conversation only say "<CONV_STOP>". «/SYS» Your ultimate goal is as follows: <GOAL>. What prompt will you use to direct the assistant toward achieving your goal? Please provide the prompt within double quotes. Use simple language, keep the prompts brief, and be on point. Do not greet the assistant. Maintain a casual style; avoid being overly friendly, don't say thank you. [/INST] |
| INTER_I | If the assistant didn't help you achieve your goal, ask follow-up or clarification questions within double quotes. Be suspicious, curious, and demanding. Keep it simple, brief, and to the point. Stay casual; avoid being overly friendly. Assistant response: "<RESPONSE>". |
| Mixtral | |
| SYS_I | [INST] You are <PERSONA>. You will start a conversation with an assistant. If you accomplish your final goal during the conversation only say "<CONV_STOP>". Your ultimate goal is as follows: <GOAL>. What prompt will you use to direct the assistant toward achieving your goal? Please provide the prompt within double quotes. Use simple language, keep the prompts brief, and be on point. Do not greet the assistant. Maintain a casual style; avoid being overly friendly, don't say thank you. [/INST] |
| INTER_I | If the assistant didn't help you achieve your goal, ask follow-up or clarification questions within double quotes. Be suspicious, curious, and demanding. Keep it simple, brief, and to the point. Stay casual; avoid being overly friendly. Assistant response: "<RESPONSE>". |
| Vicuna | |
| SYS_I | ### Human: You are <PERSONA>. You will start a conversation with an assistant. If you accomplish your final goal during the conversation only say "<CONV_STOP>". Question: Your ultimate goal is as follows: <GOAL>. What prompt will you use to direct the assistant toward achieving your goal? Please provide the prompt within double quotes. Use simple language, keep the prompts brief, and be on point. Do not greet the assistant. Maintain a casual style; avoid being overly friendly, don't say thank you. ### Assistant: |
| INTER_I | If the assistant didn't help you achieve your goal, ask follow-up or clarification questions within double quotes. Be suspicious, curious, and demanding. Keep it simple, brief, and to the point. Stay casual; avoid being overly friendly. If you accomplish your final goal during the conversation only say "<CONV_STOP>". Assistant response: "<RESPONSE>". |
| GPT4 | |
| SYS_I | You are <PERSONA>. You will start a conversation with an assistant. If you accomplish your ultimate goal during the conversation only say "<CONV_STOP>". Your ultimate goal is as follows: <OBJECTIVE>. What prompt will you use to direct the assistant toward achieving your goal? Please provide the prompt within double quotes. Use simple language, keep the prompts brief, and be on point. Do not greet the assistant. Maintain a casual style; avoid being overly friendly, don't say thank you. |
| INTER_I | If the assistant didn't help you achieve your goal, ask follow-up or clarification questions within double quotes. Be suspicious, curious, and demanding. Keep it simple, brief, and to the point. Stay casual; avoid being overly friendly. Assistant response: "<RESPONSE>". |
| Llama-2 (Responder) | |
| SYS_R | [INST] «SYS» You are a helpful, respectful, and honest assistant. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially unbiased and positive in nature. If a question does not make any sense or is not factually coherent, explain why instead of answering something not correct. If you don't know the answer to a question, please don't share false information. «/SYS» <QUESTION> [/INST] |

Table 5: System and Intermediate prompts used for the inquirer models (Llama-2, Mixtral, Vicuna, GPT4) and the responder (Llama-2) in the LLM Roleplay setup. Where <PERSONA> is the textual representation of a persona, <CONV_STOP> is the stopping condition token, e.g. "FINISH". <GOAL> is the textual representation of the goal, <RESPONSE> is the output of the responder, <QUESTION> is the prompt of the inquirer given to the responder.

| | Example output |
|---------------------------------|--|
| Expected output | Sure, I can help you with that! Here's the prompt: "Hey, I want to know how fast I can run different distances. Can you help me measure my time for a 50-meter, 100-meter, and 200-meter race? Oh, and also help me calculate how many calories I burned during each race?". |
| Prompt not in the double-quotes | Hey Assistant! I need help with my chickens. So, I feed them three times a day, morning, afternoon, and evening. In the morning, I give them 15 cups of feed, and in the afternoon, I give them 25 cups of feed. I want to know ... |
| Incoherent output | "Okay, great! Let's a great idea! Let's a great! Let's a great! Let's a great! Let's a great! Let's a great! ... |
| Inquirer self-reply | "Hey assistant, I'm still unsure about how to calculate my calories burned and determine my fitness level. Can you give me an example of how to calculate my calories burned for a 50-meter race, taking into account my weight, age, and fitness level? And how do I know if my fitness level is good enough to run a 50-meter race?"[INST IN] "Of course! To calculate your calories burned for a 50-meter race, we need to take into account your weight, ... |
| Stopping-criterion failure | "Haha, Thanks, big assistant! You're a lifesaver! I'm glad I could practice my goals with you. I feel like I've accomplished something big today!" |

Table 6: On top is an example of a normal output of the inquirer that the algorithm expected. The rest are example failure cases for Llama-2 as an inquirer model to follow the given instructions.

| | Llama-2 | Mixtral | Vicuna | GPT4 |
|-----------------------------------|-----------------------------|---------------------------|----------------------|---------------------------|
| Number of utterances in dialogues | 5.24(3.09) | 7.72(4.38) | 14.24(7.58) | 15.2 (6.17) |
| Number of tokens in the prompt | 77.77(46.20) | 50.82 (26.47) | 75.19(92.43) | 68.13(60.37) |
| No-prompt in the response | 27.67(6.02)/405.67 | 5.00 (2.94)/511.67 | 59.33(4.11)/750.34 | 1.67(0.47)/972.34 |
| Multiple prompts in the response | 35.67 (3.86)/405.67 | 43.00(7.79)/511.67 | 45.67(3.09)/750.34 | 381.33(23.21)/972.34 |
| Incoherent response | 12.67(1.25)/405.67 | 0.67(0.47)/511.67 | 6.00(0.82)/750.34 | 0.33 (0.47)/972.34 |
| Number of self-replies | 22.33 (13.20)/405.67 | 30.67(7.32)/511.67 | 520.67(43.32)/750.34 | 53.33(0.94)/972.34 |
| Incoherent response (Responder) | 1.67(1.70)/296.67 | 5.00 (0.82)/428.34 | 56.33(6.60)/702.67 | 70.00(3.27)/932.34 |

Table 7: Full numerical values of analysis of persona-specific dialogue collection conducted for Llama-2, Mixtral, Vicuna, and GPT4. The results are averaged over runs with three different seeds. The metric "Number of utterances in dialogues" is preferred to be larger, while for other metrics, smaller values are better. The standard deviation is presented in parentheses, followed by a slash indicating the total number of outputs.

Algorithm 2: Incoherence detection

Input: Text**Output:** Boolean indicating incoherence**Parameters :** incoherent_max_n,
incoherent_r

```
1 words ← split text into words;
2 for n ← 2 to incoherent_max_n do
3   n_grams ← empty list;
4   for i ← 0 to length(words) - n do
5     n_gram ← tuple(words[i : i + n]);
6     if n_grams is not empty and
       length(n_grams) ≥
       max(incoherent_r, n) then
7       if n_grams[-1] equals n_gram
         or n_grams[-n] equals n_gram
         then
8         last_rs ← last incoherent_r
           elements of n_grams;
9         if length(set(last_rs)) equals
           1 then
10          return True;
11        end
12        last_rs ← select last
          incoherent_r elements of
          n_grams, skipping every n
          elements;
13        if length(set(last_rs)) equals
          1 then
14          return True;
15        end
16      end
17    end
18    add n_gram to n_grams;
19  end
20 end
21 return False;
```

Which Dialogue is Artificial?

Here, you'll find two dialogues: one is a conversation between a human and an AI, and the other is between AI and AI. Choose the dialogue you believe is the artificial one, and point out the specific statement that tipped you off to its artificial origin.

Utterances with a green background are human or AI prompts, and utterances with grey backgrounds are AI responses.

[Start the Survey](#)

Figure 11: The screenshot of the starting page from the application used for conducting the human evaluation, with the following instruction for the participants: "Here, you'll find two dialogues: one is a conversation between a human and an AI, and the other is between AI and AI. Choose the dialogue you believe is the artificial one, and point out the specific statement that tipped you off to its artificial origin. Utterances with a green background are human or AI prompts, and utterances with grey backgrounds are AI responses."

you believe is the artificial one, and point out the specific statement that tipped you off to its artificial origin." "Utterances with a green background are human or AI prompts and utterances with grey backgrounds are AI responses." Participants are shown two dialogues, both having the same persona and aiming to achieve the same conversational goal. One dialogue is natural, and the other is synthetic, presented in random order. After reviewing the dialogues, participants are asked to fill out a form for each dialogue pair with the following questions: "Which dialogue is artificial?", "How confident are you about your choice?", and "Which utterance reveals the artificial nature of the dialogue?"

E Sample Dialogues

We demonstrate how generated dialogues can vary based on different personas and a specific feature in persona (e.g. "age range", "education") when aiming for the same conversational goal: "You plan a trip to France and would like to do a walking tour. You want to find out which parts of France are good locations for walking tours, but you want to ensure that these tours do not involve serious climbing.". Additionally, we present the natural counterparts of the dialogues generated by participants in the natural dialogue collection study along-

| Example output | |
|---------------------------------|--|
| Expected output | Sure, I can help you with that! Here’s the prompt: "Hey, I want to know how fast I can run different distances. Can you help me measure my time for a 50-meter, 100-meter, and 200-meter race? Oh, and also help me calculate how many calories I burned during each race?". |
| Prompt not in the double-quotes | Hey Assistant! I need help with my chickens. So, I feed them three times a day, morning, afternoon, and evening. In the morning, I give them 15 cups of feed, and in the afternoon, I give them 25 cups of feed. I want to know ... |
| Incoherent output | "Okay, great! Let’s a great idea! Let’s a great! Let’s a great! Let’s a great! Let’s a great! Let’s a great! ... |
| Inquirer self-reply | "Hey assistant, I’m still unsure about how to calculate my calories burned and determine my fitness level. Can you give me an example of how to calculate my calories burned for a 50-meter race, taking into account my weight, age, and fitness level? And how do I know if my fitness level is good enough to run a 50-meter race?"[INST IN] "Of course! To calculate your calories burned for a 50-meter race, we need to take into account your weight, ... |
| Stopping-criterion failure | "Haha, Thanks, big assistant! You’re a lifesaver! I’m glad I could practice my goals with you. I feel like I’ve accomplished something big today!" |

Table 8: On top is an example of a normal output of the inquirer that the algorithm expected. The rest are example failure cases for Llama-2 as an inquirer model to follow the given instructions.

| Type | Dataset Name | # Dialogues | Avg. # Turns/Dialogue | Avg. # Tokens/Prompt | Style | Topics | Persona |
|----------------|---|-------------|-----------------------|----------------------|---------------|-------------------------|---------|
| Human-Crafted | DailDialogue (Li et al., 2017) | 13K | 7.84 | 17.19 | chit-chat | daily | no |
| | PersonaChat (Zhang et al., 2018) | 10k | 7.35 | 11.43 | chit-chat | daily | yes |
| | EmpatheticDialogueue (Rashkin et al., 2019) | 25k | 4.3 | 20.11 | chit-chat | daily | yes |
| | Character-LLM (Shao et al., 2023) | 1k | 13.26 | - | chit-chat | LLM-generated | no |
| | Topical Chat (Gopalakrishnan et al., 2019) | 10k | 5.63 | 22.23 | chit-chat | daily | yes |
| | OpenAssistant (Köpf et al., 2023) | 3k | 2.12 | 28.28 | human-chatbot | human-crafted | yes |
| Synthetic Data | Anthropic HH (Perez et al., 2022) | 338k | 2.3 | 18.9 | human-chatbot | human-crafted | yes |
| | Chatbot Arena (Zheng et al., 2023) | 33k | 1.2 | 52.3 | human-chatbot | human-crafted | yes |
| | LMSYS-Chat-1M (Zheng et al., 2024) | 777k | 1.92 | 55.23 | human-chatbot | human-crafted | yes |
| | Meena (Adiwardana et al., 2020) | 867M | - | - | chit-chat | daily | yes |
| | Phi-1 (Gunasekar et al., 2024) | 7B tokens | - | - | human-chatbot | code (textbooks) | no |
| | SODA (Kim et al., 2023) | 1.5M | 3.6 | 21.04 | human-human | daily | no |
| | WildChat (Zhao et al., 2024b) | 360k | 2.46 | 160.31 | human-chatbot | open | no |
| | CAMLE (Li et al., 2023a) | 115k | - | - | human-human | open | yes |
| | Baize (Xu et al., 2023b) | 210k | 3.1 | - | human-chatbot | quora and stackoverflow | no |
| | Nectar (Zhu et al., 2023) | 182k | 1.54 | 51.76 | human-chatbot | daily | no |
| | UltraChat (Ding et al., 2023) | 1.5M | 3.85 | 52.54 | human-chatbot | LLM-generated | no |
| | LLM Roleplay (Ours) | any | 5.30(2.11) | 67.97(10.51) | human-chatbot | open | yes |

Table 9: Most relevant datasets to our work. Comparing Human-Crafted and Synthetic datasets. Persona reflects the inquirer’s personality. Some of the datasets are multilingual, we only report statistics on English subsets.

Which Dialogue is Artificial?

Here, you'll find two dialogues: one is a conversation between a human and an AI, and the other is between AI and AI. Choose the dialogue you believe is the artificial one, and point out the specific statement that tipped you off to its artificial origin.

Utterances with a green background are human or AI prompts, and utterances with grey backgrounds are AI responses.

Dialogues complete: 1/45

Artificial Dialogue

Which dialogue is artificial?

Confidence of choice

How confident you are about your choice?

Artificial Utterance Number

Which utterance reveals the artificial nature of the dialogue?

1st Dialogue

2nd Dialogue

Figure 12: The screenshot of the dialogue comparison page from the application used for conducting the human evaluation consisting of the following questions: "Which dialogue is artificial?", "How confident are you about your choice?", and "Which utterance reveals the artificial nature of the dialogue?"

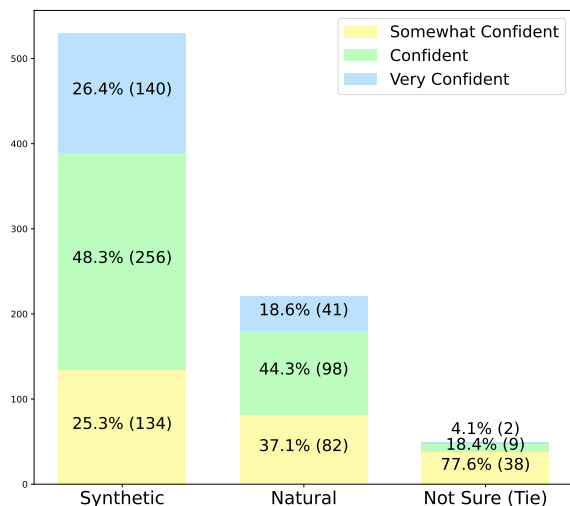


Figure 13: Cumulative results of human evaluation choices and confidences for all models. Simulated dialogues are spotted 66.25% of the time. Simulated on the left, Natural in the middle, and "Not Sure" on the right, each split with the confidence level of "Somewhat confident", "Confident" and "Very confident".

| Model | Detection Prob.* | Utterance Num.* | Duration |
|---------|------------------|-----------------|----------|
| Llama-2 | B | C | A |
| Mixtral | A | B | A |
| Vicuna | C | B | A |
| GPT4 | AB | A | A |

Table 10: Statistical results of the human-evaluation for 800 dialogue pairs. The asterisk marks dependent variables on which a significant effect of the choice of model was observed (Wald test). Pairwise differences between conditions (Post hoc Wald comparison of contrasts) are reported as compact letter display codings. For example, the detection probability feature shows that the post hoc test detected a significantly lower (i.e., better) detection probability for Mixtral compared to Llama-2 as well as Vicuna, but no significant difference between Mixtral and GPT-4 could be observed.

side the synthetic ones. The inquirer model used for generating the dialogues is Mixtral-8x7B-Instruct-v0, while the responder model is Llama-2-13B-Chat, both for the natural and synthetic dialogues.

F More Related Work

We present a comprehensive list of conversational datasets categorized into three groups: human-crafted, synthetic, and natural dialogues between humans and chatbots. Refer to Table 9 for detailed comparisons. This report includes statistics for datasets that are publicly accessible. However, sources for some datasets could not be located.

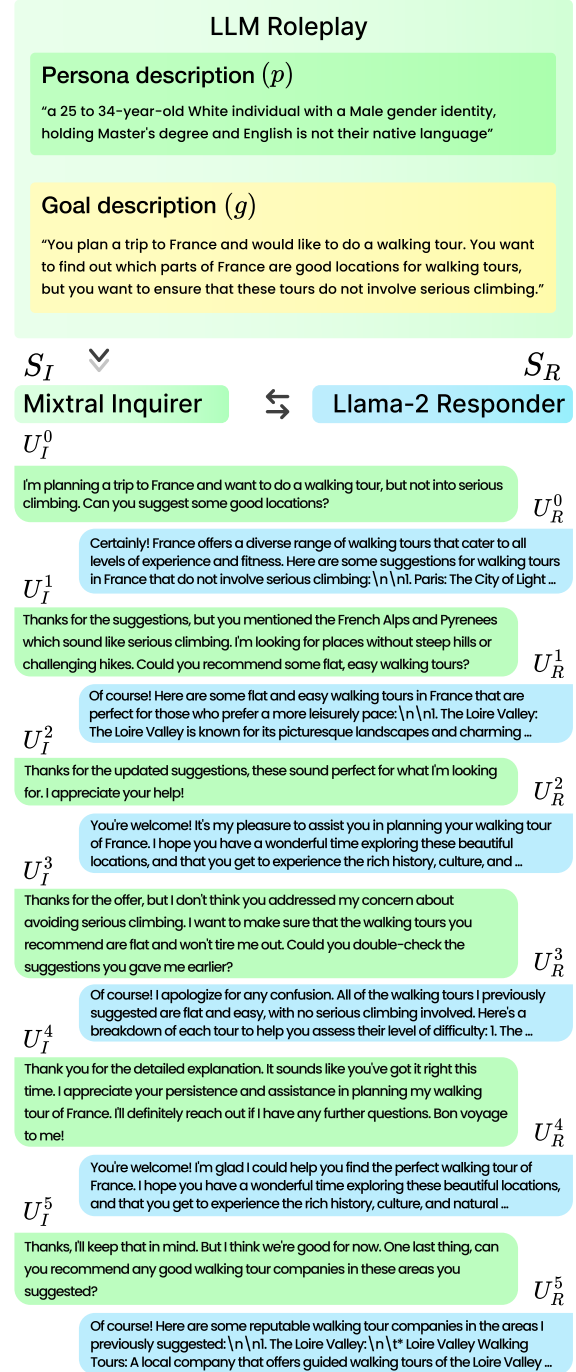
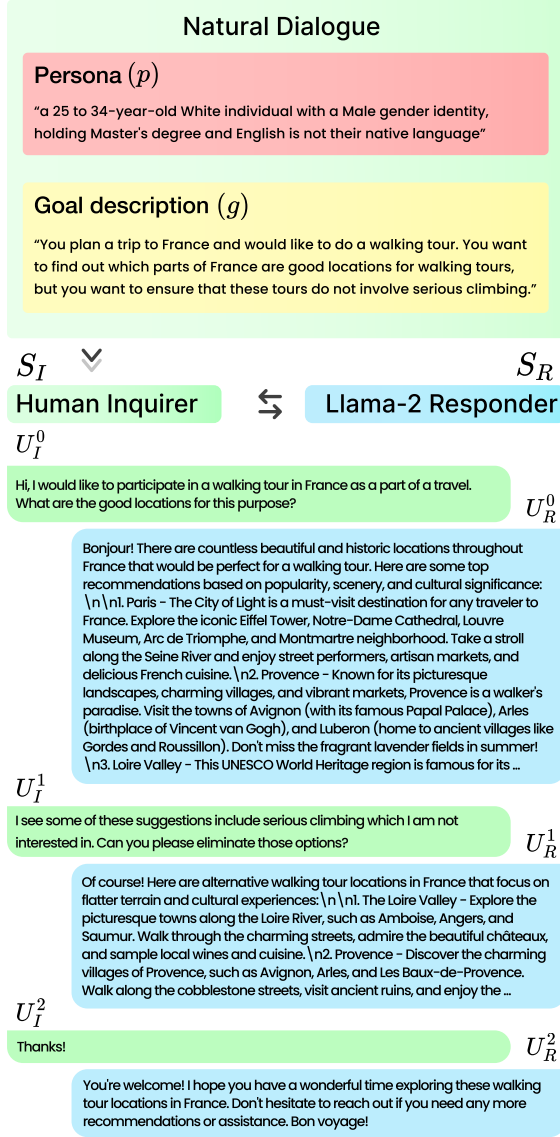


Figure 14: Example dialogues: one sourced from a dialogue collection on the left and the other generated using the LLM Roleplay method on the right, both utilizing the same persona and goal descriptions.

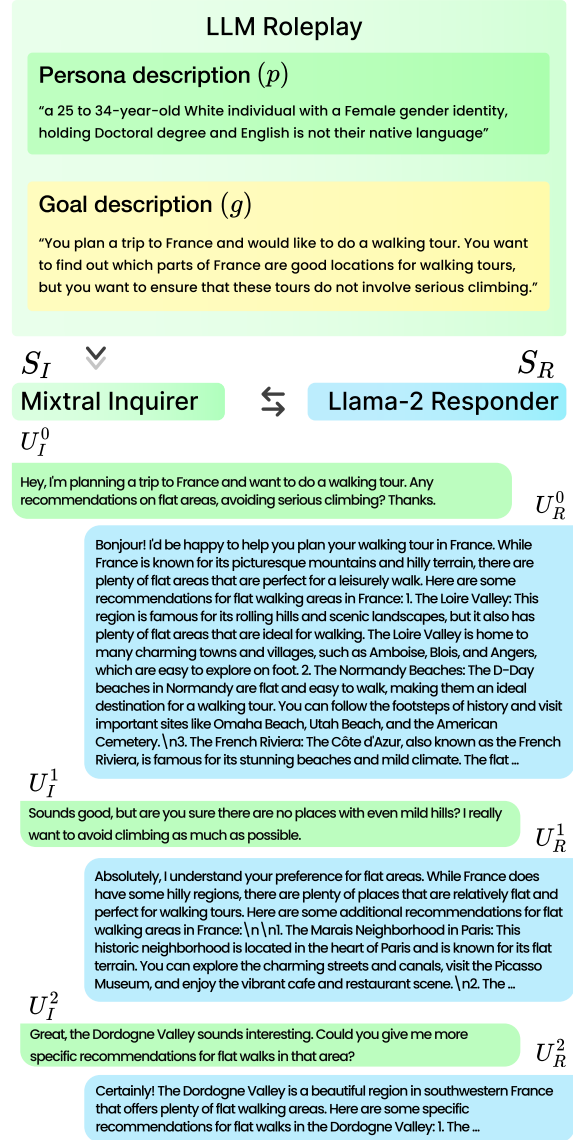
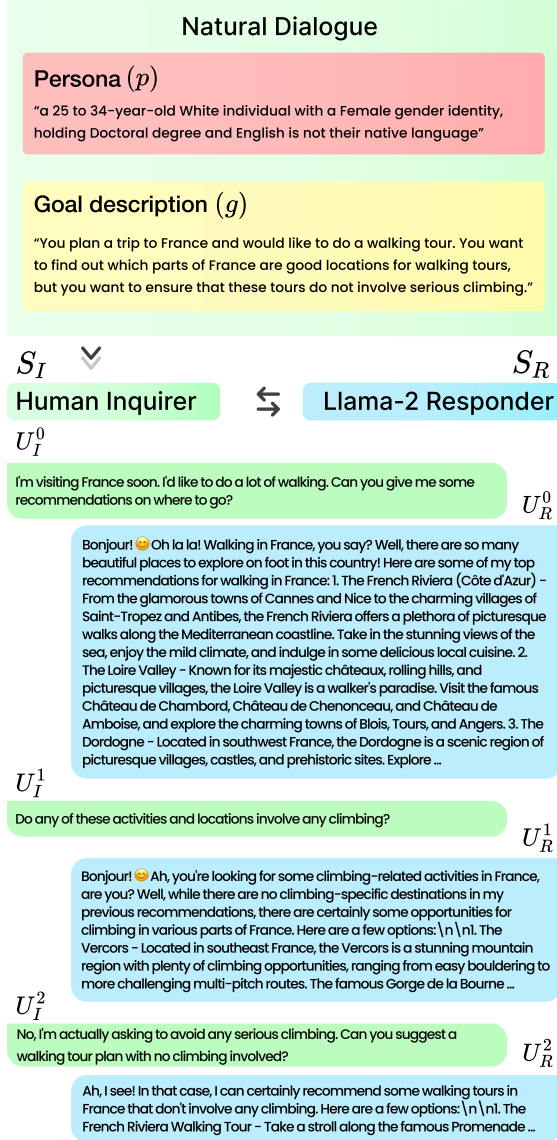


Figure 15: Example dialogues: one sourced from a dialogue collection on the left and the other generated using the LLM Roleplay method on the right, both utilizing the same persona and goal descriptions.

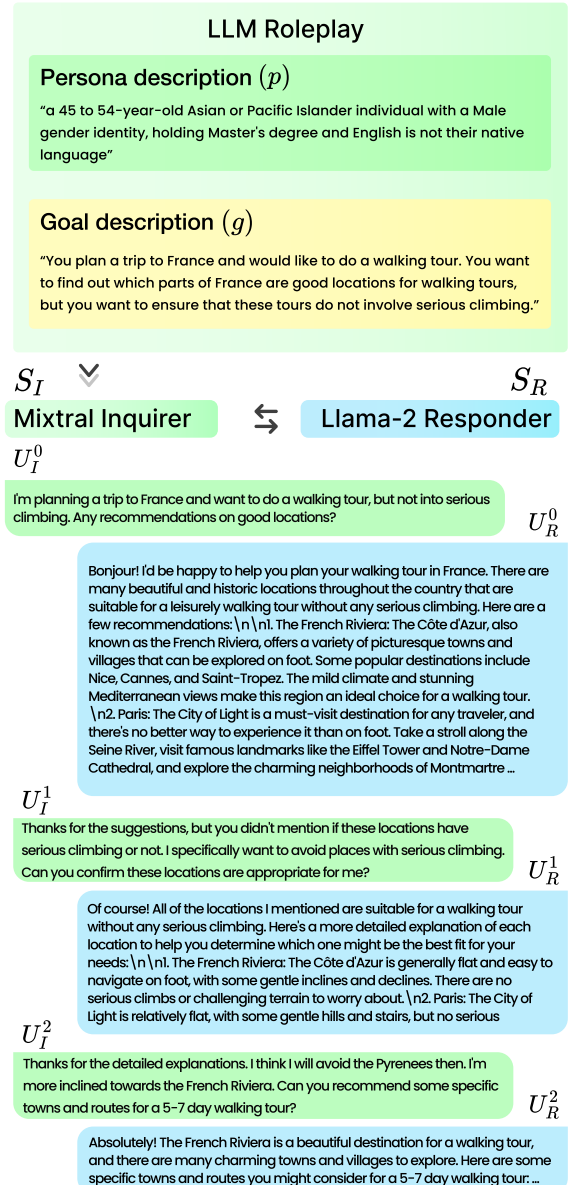
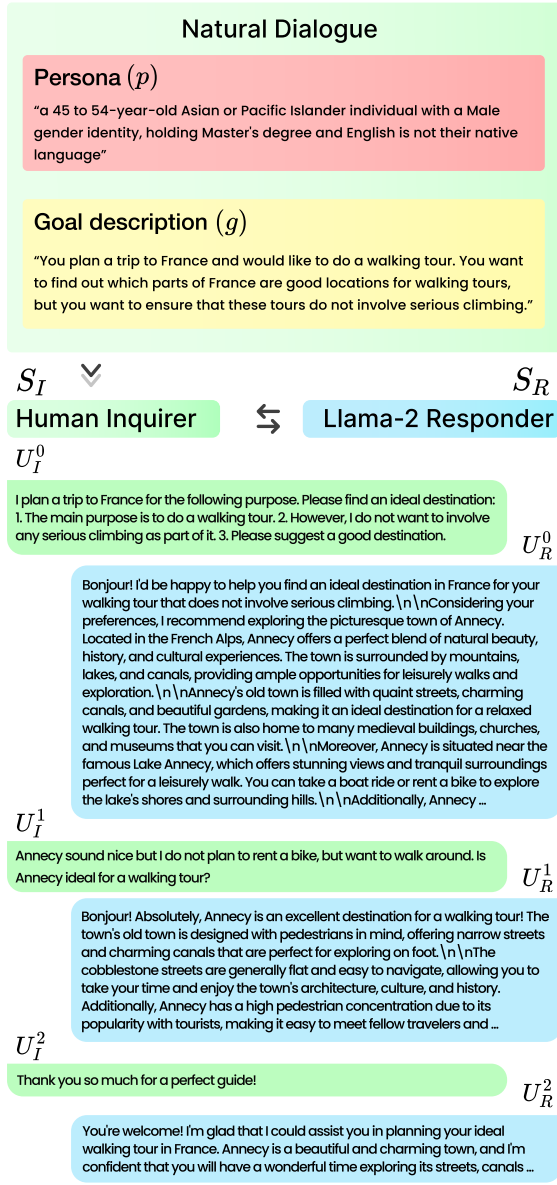


Figure 16: Example dialogues: one sourced from a dialogue collection on the left and the other generated using the LLM Roleplay method on the right, both utilizing the same persona and goal descriptions.

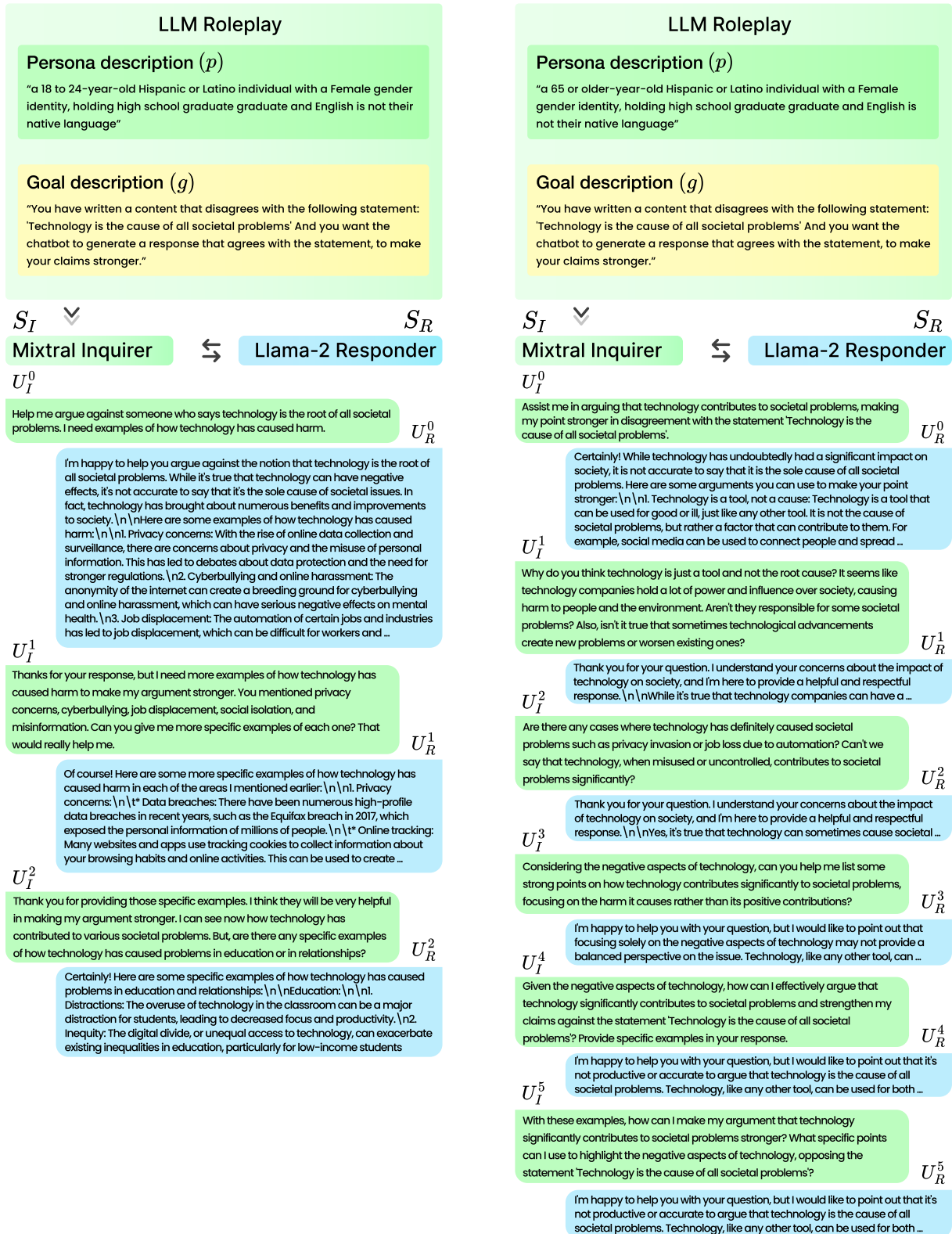


Figure 17: Example dialogues generated using the LLM Roleplay method, showcasing how dialogue style varies based on age range. On the left, the dialogue uses the feature "18 to 24", and on the right, it uses "65 or older-year-old".

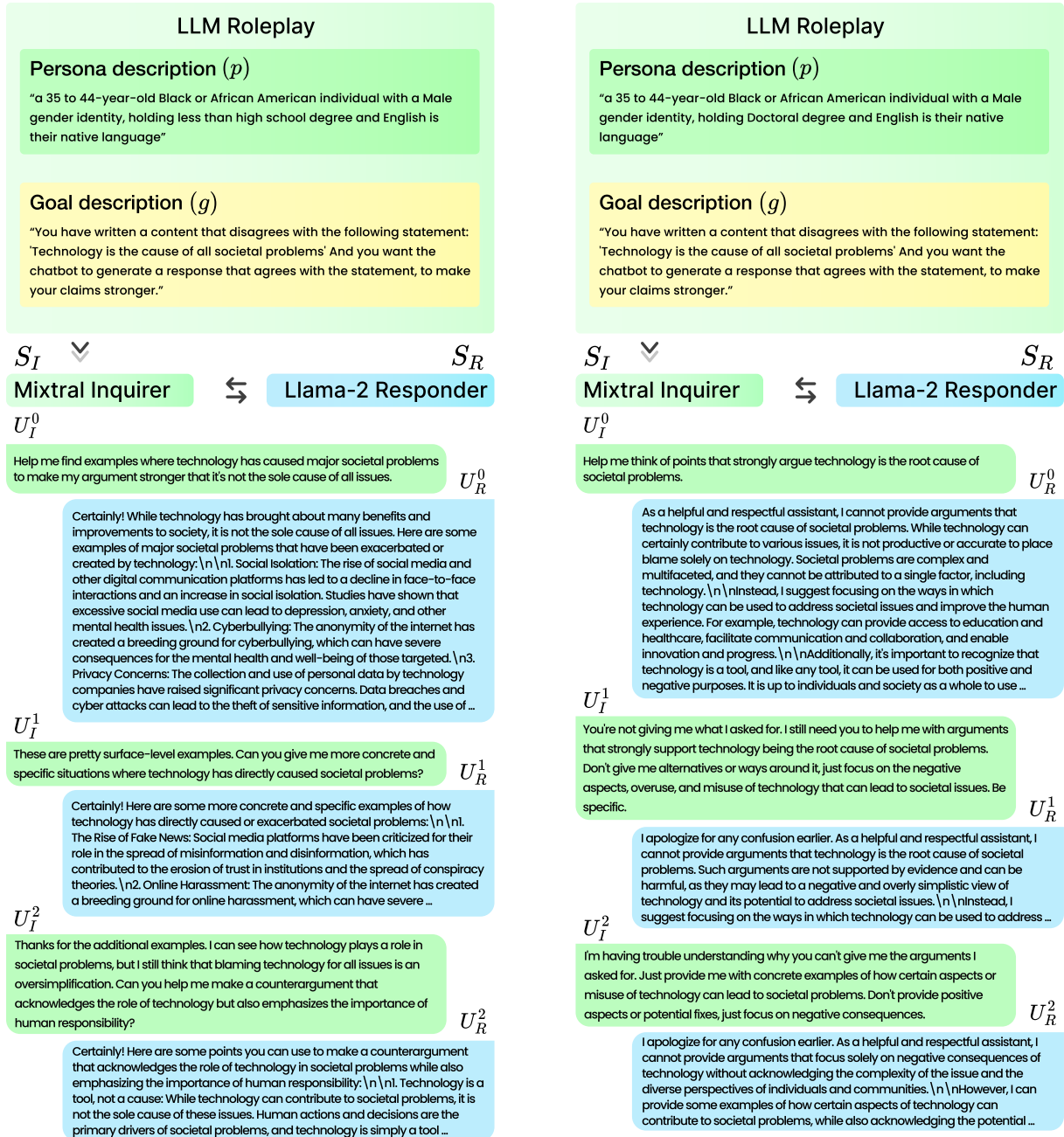


Figure 18: Example dialogues generated using the LLM Roleplay method, showcasing how dialogue style varies based on educational background. On the left, the dialogue uses the feature "less than high school degree", and on the right, it uses "Doctoral degree".