# Know You First and Be You Better: Modeling Human-Like User Simulators via Implicit Profiles

Kuang Wang<sup>1</sup>, Xianfei Li<sup>1</sup>, Shenghao Yang<sup>1</sup>, Li Zhou<sup>1</sup>, Feng Jiang<sup>2,1\*</sup>, Haizhou Li<sup>1,3</sup>

<sup>1</sup>SRIBD, School of Data Science, The Chinese University of Hong Kong, Shenzhen, Guangdong <sup>2</sup>Artificial Intelligence Research Institute, Shenzhen University of Advanced Technology <sup>3</sup>Department of ECE, National University of Singapore kuangwang@link.cuhk.edu.cn, jiangfeng@suat-sz.edu.cn

#### Abstract

User simulators are crucial for replicating human interactions with dialogue systems, supporting both collaborative training and automatic evaluation, especially for large language models (LLMs). However, current role-playing methods face challenges such as a lack of utterance-level authenticity and user-level diversity, often hindered by role confusion and dependence on predefined profiles of wellknown figures. In contrast, direct simulation focuses solely on text, neglecting implicit user traits like personality and conversation-level consistency. To address these issues, we introduce the User Simulator with Implicit Profiles (USP), a framework that infers implicit user profiles from human-machine interactions to simulate personalized and realistic dialogues. We first develop an LLM-driven extractor with a comprehensive profile schema, then refine the simulation using conditional supervised finetuning and reinforcement learning with cycle consistency, optimizing at both the utterance and conversation levels. Finally, a diverse profile sampler captures the distribution of realworld user profiles. Experimental results show that USP outperforms strong baselines in terms of authenticity and diversity while maintaining comparable consistency. Additionally, using USP to evaluate LLM on dynamic multiturn aligns well with mainstream benchmarks, demonstrating its effectiveness in real-world applications. We open-source related resources in https://github.com/wangkevin02/USP.

## 1 Introduction

The user simulator is designed as a proxy for real users in interactions with large language models (LLMs). It can simulate a realistic user by generating the target user's behavior or utterances



Figure 1: Examples of different user simulators in multiturn human-LLM interactions. "OF" and "SC" represent objective facts and subjective characteristics, respectively. highlights inconsistencies with the target profile, while indicates inauthentic user imitation.

based on the specified characteristics, enabling dynamic multi-turn interactions with LLMs (Wan et al., 2022) and scene reproduction (Wang et al., 2024b). As a result, it becomes an effective alternative (Liu et al., 2023; Ferreira et al., 2024) in scenarios where real-world human-computer interaction data is difficult to obtain, especially in domains with privacy and ethical concerns, such as medical consultations (Valizadeh and Parde, 2022). It also helps Simulation-to-Reality (Sim2Real) applications, such as tutorial strategies, election simulations, and public opinion research (Liu et al., 2024; Zhang et al., 2024b; Chuang et al., 2024).

Recent advances in LLMs have spurred the de-

<sup>&</sup>lt;sup>\*</sup>Feng Jiang is the corresponding author and Shenzhen University of Advanced Technology is the corresponding affiliation.

velopment of user simulators, improving their naturalness and utility (Deng et al., 2024; Zhang et al., 2024a). Mainstream LLM-based role-playing methods (Moon et al., 2024) use predefined profiles to mimic diverse user traits. However, as LLMs are typically trained to be universally polite and helpful (Lu et al., 2024), they often lack utterancelevel authenticity and struggle with role confusion between user simulation and their inherent assistant nature (Xu et al., 2023a), as shown in Figure 1. Models like PlatoLM (Kong et al., 2024) and Parrot (Sun et al., 2024) address this by training on real human-LLM conversation datasets, but they focus only on text utterances, lacking dialogue control. This limits conversation-level consistency and diverse simulation without seed context. Additionally, they both fail to capture the authentic distribution of user-level diversity, crucial for analyzing group behavior.

To address the issues above, we believe that a user simulator knows users' intrinsic characters hidden in their conversations first, and then can provide a better simulation. Therefore, we treat user simulation as a dialogue reconstruction task and propose a novel framework named the User Simulator with Implicit Profile (USP). It is decomposed into implicit profile extraction to capture the user's underlying characteristics from the target user dialogue and conditional generation based on the profile.

In this framework, we first propose an LLMdriven profile extractor to extract implicit profiles from user conversations with a well-designed profile schema. Inspired by interpersonal interaction theory(Kruglanski and Higgins, 2013), our profile schema contains two dimensions (objective facts (OF) and subjective characteristics (SC)) with a dozen attributes to describe the user comprehensively. Different from existing works (Cheng et al., 2024b; Tu et al., 2024) using attributes as profiles, we further polish the profile attributes into natural, descriptive profiles to ensure generalization.

Then, we integrate the extracted user profiles into the user simulator through two-stage training: (1) conditional supervised fine-tuning with user profiles for utterance-level simulation, and (2) reinforcement learning with cycle consistency to align reflected profiles from simulated dialogues with given profiles for conversation-level simulation. We also implement a diverse profile sampler to capture authentic user distributions.

Our experiments demonstrate that USP im-

proves semantic and stylistic similarity in reconstructed multi-turn dialogues by approximately 34% and 43% compared to the leading baseline, with reconstruction errors reduced by half, showcasing enhanced authenticity and diversity. It achieves dialogue profile consistency comparable to GPT-40 (User w/ Profile), improving multi-turn consistency by 14% while matching single-turn performance. Additionally, USP-based multi-turn dynamic evaluation of LLMs for downstream tasks aligns closely with established benchmarks, enabling finer-grained assessment of LLM performance across diverse user groups. Our key contributions are outlined below:

- We propose a novel approach for constructing user simulators using implicit user profiles embedded in human-LLM conversations.
- We propose a framework that infers implicit user profiles, further enhanced by conditional fine-tuning and reinforcement learning with cycle consistency, to improve simulation at both the utterance and conversation levels.
- Experiments show that USP outperforms baselines in authenticity and diversity, maintains comparable consistency, and enables effective multi-turn dynamic evaluation of LLMs.

#### 2 Related Works

#### 2.1 General User Simulator

Early user simulators focused on limited action prediction using agenda-based (Schatzmann et al., 2007; Schatzmann and Young, 2009) and modelbased methods (Asri et al., 2016; Kreyssig et al., 2018), constrained by early natural language generation capabilities—for instance, generating synthetic binary preferences in conversational recommendation systems (Christakopoulou et al., 2016).

Recent advancements in LLMs enabled more sophisticated simulations of realistic conversations, offering significantly enhanced natural language flexibility. These advances include the use of LLMs for self-chat (Xu et al., 2023b) and dual LLM architectures, where separate models roleplay user and assistant based on seed conversations (Ding et al., 2023). Following these innovations, other trained user simulators, such as PlatoLM (Kong et al., 2024) and Parrot (Sun et al., 2024), learn human discourse patterns directly from human-LLM interactions in conversations.

#### 2.2 Persona-based User Simulator

Since general user simulators often struggle to capture the full spectrum of diverse user needs, it leads a growing interest in persona-based personalization to improve both controllability and diversity in simulations (Takanobu et al., 2020). Some researchers attempt to leverage goal generators (Takanobu et al., 2020) to create diverse user goals or retrieval-based personas derived from historical data (Shi et al., 2019) to guide user simulators in task-oriented dialogue (ToD) systems.

With the rise of LLMs and their strong zero-shot role-playing capabilities (Njifenjou et al., 2024), prompt-driven user simulation has become the dominant paradigm. LLMs have been used to simulate users with predefined profiles (Chuang et al., 2024), model diverse personalities and needs in ToD (Zhang et al., 2024a), and capture user preferences in conversational recommendation (Yoon et al., 2024). Our method follows this line of work, with a focus on addressing authenticity, consistency, and diversity, which remain underexplored in related studies.

#### 3 Task Definition

We formulate user simulation as a dialogue refactoring task to replicate multiturn user behavior in a target dialogue  $d_i = \{(u_{i,1}, a_{i,1}), \dots, (u_{i,n}, a_{i,n})\}$ , where  $u_{i,j}$  is the *j*-th user utterance and  $a_{i,j}$  is the corresponding response answer. Our goal is to achieve high utterance-level and dialogue-level fidelity. Formally, we minimize utterance-level distance  $D_{\text{utt}}(u_{i,j}, u'_{i,j})$  and dialogue-level distance  $D_{\text{utt}}(u_{i,j}, u'_{i,j})$  and dialogue-level distance  $D_{\text{dia}}(d_i, d'_i)$ , where  $u'_{i,j}$  is the simulated utterance and  $d'_i$  is the simulated dialogue.

Direct simulation struggles to capture personalized traits. Recent studies (Deng et al., 2024; Kong et al., 2025) demonstrate that role-playing with specific user profiles  $(p_i)$  effectively enables diverse user simulations. However, unlike well-known figures, user profiles in real-world conversations are often implicit and challenging to derive (Wang et al., 2024a).

To address this, we reformulate the task by extracting the implicit user profile from the dialogue using a profile extractor  $P_{\text{extractor}}$ , then reconstructing the target dialogue as Eq. 1.

Category	Dimension	Attributes
Objective Facts	Scene-Consistent Attributes	Age, Gender, Location, Occupation, Education, Family Relationship, Routines/Habits, Social Relationships, Other Experiences
	Scene-Related Attributes	Goals/Plans, Task Details
Subjective	Intrinsic	Big Five Personality Traits,
Characteristics	Characteristics	Language Styles

Table 1: The Designed User Profile Schema.

$$\min_{\substack{d'_i \sim P(\cdot|p_i,\pi_\theta) \\ u'_{i,j} \sim P(\cdot|c_{i,j},p_i,\pi_\theta)}} \left[ D_{\text{utt}}(u_{i,j}, u'_{i,j}) + \alpha D_{\text{dia}}(d_i, d'_i) \right]$$
(1)

where  $p_i = P_{\text{extractor}}(d_i)$ ,  $\pi_{\theta}$  represents the learnable parameters of the language model, and  $c_{i,j} = \{(u_{i,1}, a_{i,1}), \dots, (u_{i,j-1}, a_{i,j-1})\}$  denotes the ground-truth context up to the *j*-th turn. The hyperparameter  $\alpha$  balances the utterance-level and dialogue-level distances.

## 4 Modeling User Simulator with Implicit Profiles

We propose the User Simulator with Implicit Profiles (USP) framework, shown in Figure 2, to minimize the objective in Eq. 1 across four stages. First, we build a user profile extractor with a tailored schema (Section 4.1). Then, we optimize utterance-level authenticity using conditional SFT (Section 4.2), ensure conversation-level consistency via Reinforcement Learning with Cycle Consistency (RLCC) (Section 4.4), and achieve userlevel diversity through corpus distribution fitting (Section 4.3).

#### 4.1 User Profile Construction

#### 4.1.1 User Profile Schema

We believe that the user profile should reveal user characteristics from two aspects: explicit personal information and implicit communication styles. Therefore, inspired by interpersonal interaction theory (Kruglanski and Higgins, 2013), we design a user profile schema containing objective facts(OF) and subjective characteristics(SC) to represent them, as shown in Table 1.

The OF focuses on common topics in human conversation (Cheng et al., 2024b; Dunbar et al., 1997), including scene-consistent attributes (such as age, gender, and location) and scene-related attributes (such as the goal and task details). SC



Figure 2: Overview of our proposed User Simulator with implicit Profile(USP) framework.

encompasses external personality dimensions, reflected in language style (Wang et al., 2024a), and internal personality dimensions, captured by the Big-Five Traits (Gosling et al., 2003).

Unlike prior work (Cheng et al., 2024b; Tu et al., 2024) that relies on discrete attributes for user profiles, we further reformulate these attributes into coherent narrative descriptions to enhance generalization and flexibility.

#### 4.1.2 User Profile Extractor

To obtain such a user profile, we design an LLMdriven user profile extractor that extracts the implicit user profile from the human-LLM conversation. The extractor first leverages advanced LLM (such as GPT-40) to extract the user character attributes mentioned above with a well-designed prompt. Then, the extractor collects the valid attributes together and polishes them into natural language descriptions. Further prompt details regarding the extractor can be found in Appendix A.2.

#### 4.1.3 Profile Quality Verification

Existing role-playing methods rely on predefined profiles and dialogues from separate sources, either extracted from novel segments or generated by LLMs (Zhou et al., 2024; Chen et al., 2023; Shao et al., 2023) without verifying alignment (Wang et al., 2024a; Cheng et al., 2024b). This overlooks the correlation between profiles and dialogues, potentially hindering simulation performance (Yu et al., 2024). To provide an automatic metric for evaluating this, we propose Dialogue Profile Consistency (DPC), which frames consistency as a retrieval task (Jandaghi et al., 2024). DPC employs an F1-score approach, assessing consistency through atomic fact verification by measuring both precision (DP.P) and recall (DP.R).

We first introduce Factual Consistency (Fact.Con), an adaptation of FactScore (Min et al., 2023) tailored for dialogue scenarios, as defined in Eq. 2. Given a target T, we assess its consistency with the source by decomposing T into atomic facts  $af_k$  using an atomic fact generator (afg). We then compute the natural language inference (NLI) score for each atomic fact with respect to the source S.

Fact.Con(S,T) = 
$$\frac{1}{|af_k|} \sum_{af_k \in afg(T)} NLI(S, af_k)$$
(2)

where  $NLI(\cdot, \cdot)$  denotes the NLI model, implemented using prompt-based GPT-40.

We then define  $DP.P_i = Fact.Con(d_i, p_i)$ , which measure the accuracy of profile description information and DP.R<sub>i</sub> = Fact.Con $(p_i, d_i)$  to assess the profile's coverage of the dialogue. The DPC is their harmonic mean. When dialogue  $d_i$  serves as the target T, each user utterance  $u_{i,j}$  is treated directly as an atomic fact  $af_k$ . Conversely, when the profile serves as the target T, we utilize afg followed (Min et al., 2023) to decompose it into atomic facts.

Additionally, we use a Validation Score (Val.Score) to evaluate SC description quality based on the dialogue, rated on a 1–5 scale using GPT-40(prompts detailed in Appendix D).

#### 4.2 Conditional Supervised Fine-Tuning

To empower the LLM with the general capability to simulate diverse users at the utterance level, we utilize conditional supervised fine-tuning based on user profiles. It enables the LLM to learn the conditional generation mapping based on both the extracted profile  $p_i$  and context  $c_{i,j}$ . As a subtle misalignment between the core objectives of the user simulator and the response model, the SFT language modeling loss focuses on optimizing user utterance as shown in Eq. 3.

$$\min_{\pi_{\theta}} \sum_{i,j,k} -\log P(u_{i,j,k} | u_{i,j,< k}, c_{i,j}, p_i, \pi_{\theta}) \quad (3)$$

where  $u_{i,j,k}$  denotes the k-th token of the  $u_{i,j}$ .

#### 4.3 Diverse Profile Sampling

We propose Diverse Profile Sampling to generate naturalistic user profiles that reflect real-world characteristic distributions. Our method first embeds constructed profiles into a semantic space using SimCSE (Gao et al., 2021), followed by dimensionality reduction via UMAP (McInnes et al., 2018). We then apply Gaussian Kernel Density Estimation (GKDE) to fit the underlying distribution, allowing probabilistic sampling of realistic profiles for downstream tasks such as majority representation. To further enhance diversity, we synthesize virtual profiles by combining OF and SC descriptions from nearest neighbors, producing novel yet plausible profile variants.

## 4.4 Reinforcement Learning with Cycle Consistency

The conditional SFT stage enables the user simulator to generate human-like utterances focused on forward consistency, producing precise responses aligned with the target profile. However, it does not ensure full reflection of the profile, i.e., backward consistency or profile recall. To overcome this and improve conversation-level consistency, we introduce Reinforcement Learning with Cycle Consistency (RLCC), which enhances alignment between the user simulator's actual behavior reflected in simulated dialogues and the target behavior defined by the profile, ensuring a closer match to the intended target profile.

In this stage, the user simulator  $u'_i$  interacts with a response LLM based on a target profile  $p_i$ , sampled via Diverse Profile Sampling, to generate a simulated dialogue  $d'_i$ . The dialogue ends when it reaches the maximum context length or a predefined turn limit (set to 10). The reflected profile  $p'_i$ is then extracted from  $d'_i$  using the profile generator. Our goal is to maximize the semantic similarity between the target profile  $p_i$  and the reflected profile  $p'_i$ , both in objective facts and subjective characteristics, as defined in Eq. 4.

$$\max_{\pi_{\theta}} \mathbb{E}_{p_i \sim D, d'_i \sim \pi_{\theta}(p_i)} \Big[ \operatorname{Sim}(p_i, P_{extractor}(d'_i)) \Big]$$
(4)

where  $Sim(\cdot, \cdot)$  is a similarity model (Sim-CSE (Gao et al., 2021)), and *D* denotes the virtual profiles dataset from the sampler in Section 4.3. The dialogue-level reward is uniformly attributed to each user utterance, defined as  $r_{i,j}^{cc} =$  $Sim(p_i, P_{extractor}(d'_i)).$ 

To prevent reward hacking, the AI detection reward is included as an auxiliary component. The final reward, defined in Eq. 5, is utilized to optimize profile recall via Proximal Policy Optimization (PPO)(Schulman et al., 2017).

$$r_{i,j} = \lambda r_{i,j}^{cc} + (1 - \lambda) r_{i,j}^{ai\_detect}$$
(5)

where  $r_{i,j}^{ai\_detect} = AI\_detect(u'_{i,j})$ , and  $\lambda = 0.8$ prioritizes cycle consistency. The AI\\_detect refers to a binary AI detection model (Yang et al., 2024) that predicts the probability of an utterance being AI-generated. Both the AI detection model and profile generator are fine-tuned on our training dataset, with details provided in Appendix B.1.

#### **5** Experiments

We evaluate user simulators on authenticity, consistency, and multi-turn continuity at both utterance and conversation levels, while measuring diversity by comparing the dialogue distributions of simulated and real users.

#### 5.1 Datasets

We select the popular LMSYS-Chat-1M (Zheng et al., 2023), which contains one million human-LLM conversations. Following prior work (Kong et al., 2024), we filter out non-English, toxic, and redundant samples, resulting in 94,874 conversations (87,882 for training, 4,626 for validation, and 2,366 for testing). Each conversation is then annotated with user profiles using the GPT-4obased extractor described in Section 4.1, forming the LMSYS-USP dataset. Detailed preprocessing steps are provided in Appendix A.1.

We use DPC and Val.Score to automatically evaluate the quality of extracted user profiles on 100 randomly selected samples from the LMSYS-USP test set, as well as on 100 samples each from Persona-Chat (Zhang et al., 2018) and ConvAI2<sup>1</sup> (Dinan et al., 2019), which include manually annotated predefined profiles. As shown in Table 2, the extracted profiles achieve over 84% DPC, with even distill-llama3 results comparable to GPT-40, demonstrating the effectiveness of our annotation method. Manual evaluation further confirms profile quality, with average scores exceeding 4 out of 5 (see Appendix B.3 for details).

#### 5.2 Configurations

We train USP based on LLaMA-3-8B-Base (AI@Meta, 2024) model. The conditional SFT is conducted on the training dataset using 4 A100 40GB GPUs, with full fine-tuning over 3 epochs at a learning rate of 5e-5 and max length set to 4096, taking approximately two days. Our diverse profile sampler then randomly selects 1,000 samples from the training set for virtual user sampling, combining objective facts and subjective characteristics to generate about 1 million profiles. From these, we select the 5,000 profiles least similar to the training dataset for the RLCC phase. RLCC training uses two H20 96GB GPUs for 5 days, utilizing a KL coefficient of 0.01, a learning rate of 5e-7, and training for 1 epoch.

#### 5.3 Baseline Models

(1) User Simulator without User Profile: This includes untrained GPT-40 (User w/o Profile), which use GPT-40 to predict user utterances based solely on context, and PlatoLM (Kong et al., 2024), a baseline fully fine-tuned on our dataset

using LLaMA-3-8B-Base, representing a profileagnostic approach.

(2) User Simulator with User Profile: We employ GPT-40 (User w/ Profile) and LLaMA3 (User w/ Profile) leveraging GPT-40 and LLaMA-3-8B-Instruct (AI@Meta, 2024) as profile-conditioned role-playing agents, alongside CharacterGLM (Zhou et al., 2024), a flexible profile-based baseline, and CharacterLLM (Shao et al., 2023), designed to emulate public figures.

#### 5.4 Metrics

Authenticity: We evaluate semantic and stylistic similarity using SimCSE (Gao et al., 2021) and style embeddings (Wegmann et al., 2022), respectively, to compute  $D_{utt}(u_{i,j}, u'_{i,j})$  and  $D_{dia}(d_i, d'_i)$ . To assess stylistic consistency, we report Author Verification Accuracy (AVA) (Wegmann et al., 2022), which measures whether sentence pairs are attributed to the same author based on similarity thresholds. Dialogue-level distances are computed by concatenating all user utterances.

**Consistency:** We evaluate profile-based generation consistency using reverse metrics: r-DP.P and r-DP.R. Unlike the DPC series, which treats dialogue as ground truth to assess profile quality, these metrics measure factual alignment from the profile's perspective. Specifically, r-DP.P is defined as Fact.Con $(p_i, d'_i)$ , and r-DP.R as Fact.Con $(d'_i, p_i)$ . Their harmonic mean, r-DPC, captures overall consistency. For utterance-level analysis, we report the average DP.P. Additionally, we use Persona Coverage (P.Cover) (Song et al., 2019) for keyword match and the GPT-4o-rated Subjective Characteristic Score (SC.Score) (see Appendix D) to assess subjective trait expression performance.

**Diversity:** We measure the Absolute Difference Value (ADV), defined as the Euclidean distance between PCA-reduced embeddings of generated and target dialogues, to quantify the distributional discrepancy between simulated and target dialogues.

**Continuity:** Multi-turn dialogue continuity ability is evaluated via the early stop rate (ESR), which detects premature endings triggered by repetitive responses or repeated expressions of gratitude across three turns.

#### 5.5 Results

#### 5.5.1 Utterance-Level Evaluation

In utterance-level evaluation, we assess the quality of single-turn responses generated by different user simulators given the golden context.

<sup>&</sup>lt;sup>1</sup>We use the human-to-bot dataset from https: //huggingface.co/datasets/convai-challenge/ conv\_ai\_2

Dataset	Profile Source		SC				
		DP.P↑	Avg DP.P # Fact	DP.R↑	Avg DP.R # Fact	DPC↑	Val.Score↑
LMSYS-USP	GPT4o	86.89	25.64	82.24	3.71	84.50	4.42
LMSYS-USP	Distill-llama3	86.15	23.81	81.95	3.71	84.00	4.36
Persona Chat	GPT40	86.21	22.82	62.76	7.86	72.64	4.35
Persona Chat	Human	76.21	8.59	42.94	7.86	54.93	-
ConvAI2	GPT40	68.71	17.44	39.15	9.97	49.88	3.47
ConvAI2	Human	25.69	8.70	12.64	9.97	16.94	-

Table 2: Automated evaluation of profile quality across datasets. Avg DP.P # Fact denotes the average number of atomic facts per user profile, while Avg DP.R # Fact represents the average number of user utterances per dialogue. Note that human-annotated profiles in PersonaChat and ConvAI2 contain no subjective characteristics.

Model Type	Model		Authenticity		Consistency			
		Sem-Sim↑	Style-Sim↑	AVA↑	r-DP.P↑	P.Cover↑	SC.Score↑	
w/a Drafila	GPT-40 (User w/o Profile)	40.24	13.75	11.28	_	_	-	
w/o Ploine	PlatoLM	39.37	43.11	40.29	-	-	-	
	Character_LLM	37.54	18.88	15.03	54.77	66.62	2.43	
	Character_GLM	38.51	22.28	18.17	68.72	57.72	2.95	
w/ Profile	LLaMA3 (User w/ Profile)	39.82	14.88	13.47	82.19	72.29	3.92	
	GPT-40 (User w/ Profile)	41.66	5.74	9.87	92.73	73.34	4.71	
	USP (w/o RLCC)	54.25	46.57	43.61	71.30	71.56	3.36	
	USP	53.38	46.60	43.35	72.61	71.23	3.39	

Table 3: Utterance-level performance comparison of different user simulator.

As shown in Table 3, USP outperforms all baselines in authenticity, as measured by both semantic similarity (Sem-Sim: 53.38) and stylistic similarity (Style-Sim: 46.60). This shows the effectiveness of our implicit profile-based approach for user-LLM dialogue reconstruction, especially compared to non-profile baselines like PlatoLM. While dedicated role-playing models (e.g., GPT-40 (User w/ Profile)) achieve higher consistency scores (r-DP.P) due to direct profile keyword copying with high P.Cover (73.34), our USP strikes a better balance between authenticity and consistency, as shown by the intuitive examples in Section 6.3.

#### 5.5.2 Conversation-Level Evaluation

In the conversation-level evaluation, we assess the quality of multi-turn dialogues generated by different user simulators interacting with GPT-40, each provided with either a given profile or the first turn of a reference dialogue.

As shown in Table 4, USP outperforms baseline models in authenticity, consistency, and continuity. With the lowest ESR (10), USP ensures superior dialogue continuity, avoiding issues like repetitive generation and reciprocal appreciation loops seen in baselines. Its advantage in authenticity is especially evident in multi-turn scenarios, compared to sentence-level evaluations. In terms of consistency, USP excels with a high r-DP.R (74.38) and significantly better r-DPC (64.05), demonstrating strong conditional generation consistency. Unlike role-playing models such as GPT-40 (User w/ Profile) and LLaMA3 (User w/ Profile), which show high P.Cover and r-DP.P but lower overall profile dialogue consistency, USP demonstrates a deeper and more comprehensive understanding of user behavior, moving beyond surface-level keyword matching to deliver a more vivid user simulation.

#### 5.5.3 Human Evaluation

We randomly selected 100 samples from the test set and engaged 8 evaluators to assess conversations on authenticity and consistency. Authenticity was evaluated based on Style, Semantics, and Quality, while consistency focused on Accuracy, Completeness, and Quality. Detailed criteria are in Appendix B.3.

Table 5 shows USP's clear superiority in both authenticity and consistency. USP outperforms GPT-40 (User w/ Profile) in authenticity (74 vs. 13) and consistency (61 vs. 35). It also surpasses PlatoLM trained on the same data in authenticity, demonstrating the advantage of implicit profile modeling. The larger gap in consistency (43 vs. 30) compared to authenticity (37 vs. 31) between USP and USP (w/o RLCC) highlights RLCC's key

Model Type	Model	Continuity	Authenticity			Consistency				
niouer type		ESR↓	Sem-Sim↑	Style-Sim↑	AVA↑	r-DP.P↑	r-DP.R↑	r-DPC↑	P.Cover↑	SC.Score↑
w/a Drafila	GPT-40 (User w/o Profile)	35	48.91	14.21	10.58	_	-	-	-	-
w/o Ploine	PlatoLM	18	43.24	32.43	31.60	-	-	-	-	-
	Character_LLM	52	23.37	7.13	4.69	25.48	6.43	10.27	21.49	2.82
	Character_GLM	44	40.19	10.86	12.67	39.51	29.61	33.85	42.75	3.64
w/ Profile	LLaMA3 (User w/ Profile)	31	46.84	10.58	11.63	67.09	29.98	41.44	47.72	4.19
	GPT-40 (User w/ Profile)	32	48.87	10.15	11.26	76.59	43.72	55.66	51.02	4.56
	USP (w/o RLCC)	12	66.17	40.01	35.68	53.17	71.88	61.13	42.63	3.24
	USP	10	65.39	46.23	38.77	56.24	74.38	64.05	44.08	3.35

Table 4: Conversation-level performance comparison of different user simulators.

role in aligning profiles with dialogues.

Baseline	Metrics (% USP win/tie/loss)					
	Authenticity $(\kappa=0.548)$	Consistency $(\kappa=0.561)$				
GPT-40 (User w/ Profile)	74/13/13	61/4/35				
PlatoLM	55/12/33	-				
USP (w/o RLCC)	37/32/31	43/27/30				

Table 5: Human evaluation of USP win rates over baselines in terms of authenticity and consistency.  $\kappa$  denotes the within-group kappa coefficient. Note that PlatoLM, as a non-profile baseline, contain no consistency results.

#### 5.5.4 Diversity Sampling Evaluation



Figure 3: Cumulative distribution of ADV performance comparison across different simulators, with the red cross indicating that USP (w/o RLCC) has 60% of its samples with ADV below 5%.

Figure 3 compares ADV between target and generated dialogues across simulators and percentiles. The USP series consistently demonstrate lower ADV, even in extreme cases, with 60% of samples achieving ADV below 5% (marked by a red cross), compared to baselines(e.g. PlatoLM, GPT-40 (User w/ Profile)) at 10% or higher. This reflects the USP series' superior ability to generate dialogues that closely align with target conversations, effectively preserving the diversity distribution of user characteristics. We also show that our sampling strategy outperforms random sampling by effectively capturing diverse representatives (majority and minority) in Appendix C.1.

#### 6 Analysis

#### 6.1 Ablation Study

We evaluate the effectiveness of the polishing step in our two-stage profile construction pipeline, which converts attributes into natural language descriptions, by comparing it to a baseline (USAP w/o RLCC) that uses only attributes without polishing. As shown in the first two rows of Table 6, the polishing step enhances generalization, improving performance across most metrics of Continuity, Authenticity, and Consistency. In contrast, relying solely on attributes leads USAP (w/o RLCC) to excessively replicate profile descriptions, resulting in a high P.Cover score (50.35) due to attributes appearing directly in the dialogue.

We also assess the relative importance of RLCC's two rewards by testing different  $\lambda$  values in Equation 5, denoted as USP( $\lambda$ :  $1 - \lambda$ ). Table 6 shows that  $\lambda = 0.8$  optimally balances model capabilities and dialogue consistency. Higher  $\lambda$  (0.9) sacrifices speaking style authenticity without improving r-DPC, increasing P.Cover and indicating superficial profile matching. Conversely,  $\lambda = 0.5$  preserves authentic style but lacks sufficient consistency, resulting in stagnant performance.

# 6.2 Applications: Dynamic Multi-turn Evaluation For LLMs

One application of our simulator is addressing the gap in dynamic multi-turn evaluation of LLMs. While current automatic evaluations rely on static preset questions (Zheng et al., 2024; Bai et al., 2024), our user simulator can dynamically interact with the LLM over multiple rounds, adjusting based on response quality and given user traits.

Model Configuration	Continuity	А	Authenticity			Consistency			
	ESR↓	Sem-Sim↑	Style-Sim↑	AVA↑	r-DP.P↑	r-DP.R↑	r-DPC↑	P.Cover↑	SC.Score↑
USAP (w/o RLCC)	17	64.22	40.96	35.55	65.36	51.29	57.47	50.35	3.21
USP (w/o RLCC)	12	66.17	40.01	35.68	53.17	71.88	61.13	42.63	3.24
USP (5:5)	14	66.28	41.22	37.03	52.23	71.59	60.39	43.58	3.55
USP (8:2)	10	65.39	46.23	38.77	56.24	74.38	64.05	44.08	3.35
USP (9:1)	12	66.91	38.87	33.62	58.36	70.62	63.90	46.75	3.33

Table 6: Ablation study of our USP framework.

We generated 300 diverse user profiles using our sampler: 100 highest-probability (majority), 100 lowest-probability (minority), and 100 random synthetic (virtual) profiles. Using these profiles, USP engages in multi-turn dialogues with the LLM. Evaluation results, based on MT-Bench (Zheng et al., 2024) and presented in Table 7, show that dynamic multi-turn evaluation aligns closely with average rankings on LiveBench (White et al., 2024) and Chatbot-Arena (Chiang et al., 2024), while revealing fine-grained weaknesses of different LLMs when interacting with specific user groups. Detailed analysis is in Appendix C.3.

Model Setup	Sam	pling Str	ategy	Δυσ	Ranking in LiveBench	
Model Setup	Major	Minor	Virtual	111.6.	Chatbot-Arena	
Deepseek-v3	8.25	6.13	7.70	7.36	1	
GPT-40	7.86	6.65	7.19	7.23	3	
Claude-Sonnet	7.18	6.61	7.48	7.09	2	
4o-Mini	6.84	5.70	5.52	6.02	4	
Claude-Haiku	4.88	5.42	5.43	5.24	5	

Table 7: LLM performance across user groups.

#### 6.3 Case Study

Table 8 shows that the role-playing baseline GPT-40 (w/ Profile) often copies abstract profile traits verbatim to assert user identity. In contrast, our USP conveys these traits more naturally by transforming abstract concepts into concrete and coherent expressions. For instance, when the profile states "being a father of two," GPT-40 (w/ Profile) repeats it directly, while USP implicitly reflects this by mentioning a "son" and a "daughter" in later turns. Similarly, rather than restating "likes Italian food," USP refers to a specific dish like pasta. These examples illustrate that USP better mimics human behavior by expressing abstract, high-level traits both directly and subtly in a natural manner, likely contributing to its stronger human preference (see Table 5). Additional case studies are provided in Appendix C.2.

(*Profile*) As a father of two, your strong sense of family responsibility is evident... Your interest in vegan Italianinspired snack dishes...

(1st turn) **GPT-4o (w/ Profile):** <u>As a father of two</u> who has a keen interest in ...

(2nd turn) **GPT-40** (w/ Profile): Lately, I've been exploring new vegan Italian-inspired snack dishes...

(4th turn) USP:My son wants pasta but he doesn't like tomato sauce...

(5th turn) ...my daughter says she does not feel well... She prefers vegan food in Italian style ...

Table 8: Case study comparing outputs from USP and the GPT-40 (w/ Profile) baseline for a sample profile. **Bold** highlights keywords explicitly copied by GPT-40, while **bold italic** marks USP's implicit or fuzzy matches.

#### 7 Conclusion

In this work, we introduce the USP framework, which integrates extracted user profiles into the user simulator by conditional SFT and RLCC. Our experimental results, validated by both automatic metrics and human evaluations, show that USP significantly outperforms role-playing simulators (e.g., GPT-40 (User w/o Profile)) and direct simulation approaches (e.g., PlatoLM) in authenticity and diversity while achieving comparable consistency at both the sentence and conversation levels. Additionally, dynamic evaluations with various LLMs across diverse demographic groups highlight USP's effectiveness in real-world scenarios. Nonetheless, a gap remains compared to real human behavior, and our future work will explore finer-grained control and multimodal simulation.

## Limitations

We acknowledge the following limitations: 1) Scenario Applicability: Experiments were conducted on a single dataset, with minimal validation across others to confirm broader applicability. 2) Linguistic and Cultural Scope: Our focus on English dialogues may limit the applicability of USP to other languages and cultural contexts.

### Acknowledgments

This research is supported by the project of Shenzhen Science and Technology Research Fund (Fundamental Research Key Project, Grant No. JCYJ20220818103001002), Shenzhen Science and Technology Program (Shenzhen Key Laboratory, Grant No. ZDSYS20230626091302006), Shenzhen Stability Science Program 2023, Shenzhen Key Lab of Multi-Modal Cognitive Computing, SRIBD Innovation Fund (Grant No. K00120240006), and Program for Guangdong Introducing Innovative and Entrepreneurial Teams, Grant No. 2023ZT10X044.

#### References

AI@Meta. 2024. Llama 3 model card.

- Layla El Asri, Jing He, and Kaheer Suleman. 2016. A sequence-to-sequence model for user simulation in spoken dialogue systems. *arXiv preprint arXiv:1607.00070*.
- Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, et al. 2024. Mt-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues. *arXiv preprint arXiv:2402.14762*.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv* preprint arXiv:2004.05150.
- Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and Jia Li. 2023. Large language models meet harry potter: A dataset for aligning dialogue agents with characters. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8506–8520.
- Chuanqi Cheng, Quan Tu, Wei Wu, Shuo Shang, Cunli Mao, Zhengtao Yu, and Rui Yan. 2024a. "indialogues we learn": Towards personalized dialogue without pre-defined profiles through in-dialogue learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10408–10422.
- Yi Cheng, Wenge Liu, Kaishuai Xu, Wenjun Hou, Yi Ouyang, Chak Tou Leong, Xian Wu, and Yefeng Zheng. 2024b. Evolving to be your soulmate: Personalized dialogue agents with dynamically adapted personas. *CoRR*.
- Zihao Cheng, Li Zhou, Feng Jiang, Benyou Wang, and Haizhou Li. Beyond binary: Towards fine-grained llm-generated text detection via role recognition and involvement measurement. In *THE WEB CONFER-ENCE 2025*.

- 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.
- Konstantina Christakopoulou, Filip Radlinski, and Katja Hofmann. 2016. Towards conversational recommender systems. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 815–824.
- Yun-Shiuan Chuang, Krirk Nirunwiroj, Zach Studdiford, Agam Goyal, Vincent Frigo, Sijia Yang, Dhavan Shah, Junjie Hu, and Timothy T. Rogers. 2024. Beyond demographics: Aligning role-playing llmbased agents using human belief networks. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November* 12-16, 2024, pages 14010–14026.
- Yang Deng, Wenxuan Zhang, Wai Lam, See-Kiong Ng, and Tat-Seng Chua. 2024. Plug-and-play policy planner for large language model powered dialogue agents. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.*
- Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander H. Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W. Black, Alexander I. Rudnicky, Jason Williams, Joelle Pineau, Mikhail S. Burtsev, and Jason Weston. 2019. The second conversational intelligence challenge (convai2). *CoRR*, abs/1902.00098.
- 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.
- Robin IM Dunbar, Anna Marriott, and Neil DC Duncan. 1997. Human conversational behavior. *Human nature*, 8:231–246.
- Rafael Ferreira, David Semedo, and João Magalhães. 2024. Multi-trait user simulation with adaptive decoding for conversational task assistants. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November* 12-16, 2024, pages 16105–16130.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 6894–6910.
- Samuel D Gosling, Peter J Rentfrow, and William B Swann Jr. 2003. A very brief measure of the bigfive personality domains. *Journal of Research in personality*, 37(6):504–528.

- Pegah Jandaghi, XiangHai Sheng, Xinyi Bai, Jay Pujara, and Hakim Sidahmed. 2024. Faithful personabased conversational dataset generation with large language models. In *Findings of the Association* for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, pages 15245–15270.
- Chuyi Kong, Yaxin Fan, Xiang Wan, Feng Jiang, and Benyou Wang. 2024. Platolm: Teaching llms in multi-round dialogue via a user simulator. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 7841–7863.
- Chuyi Kong, Ziyang Luo, Hongzhan Lin, Zhiyuan Fan, Yaxin Fan, Yuxi Sun, and Jing Ma. 2025. Sharp: Unlocking interactive hallucination via stance transfer in role-playing llms.
- Florian Kreyssig, Iñigo Casanueva, Paweł Budzianowski, and Milica Gasic. 2018. Neural user simulation for corpus-based policy optimisation of spoken dialogue systems. In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*, pages 60–69.
- Arie W Kruglanski and E Tory Higgins. 2013. *Social psychology: Handbook of basic principles*. Guilford Publications.
- Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. 2024. MT-eval: A multiturn capabilities evaluation benchmark for large language models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 20153–20177.
- Yajiao Liu, Xin Jiang, Yichun Yin, Yasheng Wang, Fei Mi, Qun Liu, Xiang Wan, and Benyou Wang. 2023. One cannot stand for everyone! leveraging multiple user simulators to train task-oriented dialogue systems. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 1–21.
- Zhengyuan Liu, Stella Xin Yin, Geyu Lin, and Nancy Chen. 2024. Personality-aware student simulation for conversational intelligent tutoring systems. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024, pages 626–642.
- Keming Lu, Bowen Yu, Chang Zhou, and Jingren Zhou. 2024. Large language models are superpositions of all characters: Attaining arbitrary role-play via self-alignment. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers).
- Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei

Fang. 2024. Evaluating very long-term conversational memory of llm agents. *arXiv preprint arXiv:2402.17753*.

- Leland McInnes, John Healy, and James Melville. 2018. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023, pages 12076–12100.
- Suhong Moon, Marwa Abdulhai, Minwoo Kang, Joseph Suh, Widyadewi Soedarmadji, Eran Kohen Behar, and David M. Chan. 2024. Virtual personas for language models via an anthology of backstories. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024, pages 19864–19897.
- Ahmed Njifenjou, Virgile Sucal, Bassam Jabaian, and Fabrice Lefèvre. 2024. Role-play zero-shot prompting with large language models for open-domain human-machine conversation. *arXiv preprint arXiv:2406.18460*.
- Alex Rodriguez and Alessandro Laio. 2014. Clustering by fast search and find of density peaks. *Science*, 344(6191):1492–1496.
- Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young. 2007. Agenda-based user simulation for bootstrapping a pomdp dialogue system. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Companion Volume, Short Papers*, pages 149–152.
- Jost Schatzmann and Steve Young. 2009. The hidden agenda user simulation model. *IEEE transactions on audio, speech, and language processing*, 17(4):733–747.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-LLM: A trainable agent for roleplaying. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 13153–13187.
- Weiyan Shi, Kun Qian, Xuewei Wang, and Zhou Yu. 2019. How to build user simulators to train rl-based dialog systems. arXiv preprint arXiv:1909.01388.

- Haoyu Song, Weinan Zhang, Yiming Cui, Dong Wang, and Ting Liu. 2019. Exploiting persona information for diverse generation of conversational responses. In Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019, pages 5190–5196.
- Yuchong Sun, Che Liu, Kun Zhou, Jinwen Huang, Ruihua Song, Xin Zhao, Fuzheng Zhang, Di Zhang, and Kun Gai. 2024. Parrot: Enhancing multi-turn instruction following for large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 9729–9750.
- Ryuichi Takanobu, Runze Liang, and Minlie Huang. 2020. Multi-agent task-oriented dialog policy learning with role-aware reward decomposition. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 625–638.
- Quan Tu, Shilong Fan, Zihang Tian, Tianhao Shen, Shuo Shang, Xin Gao, and Rui Yan. 2024. CharacterEval: A Chinese benchmark for role-playing conversational agent evaluation. In *Proceedings of the* 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 11836–11850.
- Mina Valizadeh and Natalie Parde. 2022. The AI doctor is in: A survey of task-oriented dialogue systems for healthcare applications. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022,* pages 6638–6660.
- Dazhen Wan, Zheng Zhang, Qi Zhu, Lizi Liao, and Minlie Huang. 2022. A unified dialogue user simulator for few-shot data augmentation. In *Findings of the Association for Computational Linguistics: EMNLP* 2022, pages 3788–3799.
- Noah Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. 2024a. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 14743–14777.
- Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR.
- Zhenduo Wang, Zhichao Xu, Vivek Srikumar, and Qingyao Ai. 2024b. An in-depth investigation of user response simulation for conversational search. In *Proceedings of the ACM Web Conference 2024*, pages 1407–1418.

- Anna Wegmann, Marijn Schraagen, and Dong Nguyen. 2022. Same author or just same topic? towards content-independent style representations. In *Proceedings of the 7th Workshop on Representation Learning for NLP, RepL4NLP@ACL 2022, Dublin, Ireland, May 26, 2022*, pages 249–268.
- Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Siddartha Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. 2024. Livebench: A challenging, contamination-free llm benchmark.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023a. 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.
- 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.
- Lingyi Yang, Feng Jiang, Haizhou Li, et al. 2024. Is chatgpt involved in texts? measure the polish ratio to detect chatgpt-generated text. *APSIPA Transactions on Signal and Information Processing*, 13(2).
- Se-eun Yoon, Zhankui He, Jessica Maria Echterhoff, and Julian McAuley. 2024. Evaluating large language models as generative user simulators for conversational recommendation. *arXiv preprint arXiv:2403.09738*.
- Yeyong Yu, Runsheng Yu, Haojie Wei, Zhanqiu Zhang, and Quan Qian. 2024. Beyond dialogue: A profile-dialogue alignment framework towards general role-playing language model. *arXiv preprint arXiv:2408.10903*.
- 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, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers, pages 2204–2213.
- Tong Zhang, Chen Huang, Yang Deng, Hongru Liang, Jia Liu, Zujie Wen, Wenqiang Lei, and Tat-Seng Chua. 2024a. Strength lies in differences! improving strategy planning for non-collaborative dialogues via diversified user simulation. In *Proceedings of the* 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024, pages 424–444.
- Xinnong Zhang, Jiayu Lin, Libo Sun, Weihong Qi, Yihang Yang, Yue Chen, Hanjia Lyu, Xinyi Mou, Siming Chen, Jiebo Luo, et al. 2024b. Electionsim: Massive population election simulation powered by

large language model driven agents. *arXiv preprint arXiv:2410.20746*.

- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric P Xing, et al. 2023. Lmsyschat-1m: A large-scale real-world llm conversation dataset. *arXiv preprint arXiv:2309.11998*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Pei Ke, Guanqun Bi, Libiao Peng, et al. 2024. Characterglm: Customizing social characters with large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 1457–1476.

#### A Dataset Construction

#### A.1 Preprocessing

Our dataset preprocessing follows the method outlined in PlatoLM (Kong et al., 2024), which includes the removal of non-English content, filtering of toxic data, elimination of exact duplicates at the dialogue level, and segmentation of conversations into maximum-length token sequences. To maintain discourse integrity, truncated dialogues are ensured to start with the assistant's turn, preserving context consistency and coherence.

## A.2 Profile Dataset

As detailed in Section 4.1, we classify attributes into three types: scene-consistent, scene-related, and deep intrinsic characteristics. For each, we use specific prompts (Figures 11, 12, and 13), with metric definitions based on (Cheng et al., 2024b) and Big Five traits per (Gosling et al., 2003).

We then concatenate these attributes, remove invalid entries, and shuffle their order to prevent positional bias. The combined attributes are rephrased using GPT-40 with the prompt in Figure 14, producing automatically labeled profiles. The LMSYS-USP dataset averages 1,149 tokens in training, 1,295 in validation, 1,438 in testing, and 231 tokens per profile.

We also measured the frequency of each attribute value, defined as the average number of distinct values per sample, to assess attribute prevalence. Statistics for objective facts are shown in Figure 4. For subjective traits, we focused on the Big Five traits only when scores were significantly high or low, excluding moderate scores as they reflect average human behavior (Moon et al., 2024) and were omitted from the subsequent polishing step.

Attribute	High Rate (%)	Low Rate (%)
Conscientiousness	78.07	7.53
Agreeableness	6.45	14.98
Extraversion	4.08	14.15
Openness	58.77	5.30
Neuroticism	2.04	10.12

 Table 9: Summary of extracted subjective attribute statistics.



Figure 4: Frequency of values for each attribute of objective facts in the attribute extraction process.

# A.3 Resource Consumption in Implementation

Attribute extraction using the GPT-40 API costs \$0.003 per attribute type, or \$0.01 per sample for three types. For 94,000 samples, the extraction costs \$940. Rewriting attributes into profiles adds \$0.05 per sample, resulting in a total dataset construction cost of \$1,400.

#### **B** Implement Detail

#### **B.1** Trainable Model Setup

For PlatoLM, we base it on the LaMA-3-8B-Base architecture. Following (Kong et al., 2024), the system prompt is: "A chat between a curious human and an artificial intelligence assistant. The human can ask follow-up or new questions without prior context." We fine-tune using four A100 40GB GPUs for 3 epochs, taking about two days.

The AI detection model uses Longformer (Beltagy et al., 2020) trained on our dataset per (Cheng et al.). User utterances are labeled as human, and assistant utterances as AI. Training runs for 3 epochs on dual RTX 3090 GPUs, taking three days.

The profile generator is fine-tuned from LLaMA-3-8B-Instruct (AI@Meta, 2024) on our curated profile dataset, effectively distilling GPT-4o's twostage profile generation. Training uses four A100 40GB GPUs for 3 epochs, lasting two days.

#### **B.2** Train-Free Model Setup

We use two simulator types: (1) A response model, e.g., GPT-40 (User with Profile), role-playing as a user simulator, unable to initiate conversations. The user profile is embedded in the system prompt, and the initial query is generated by asking, "What will you say to start the conversation?" to obtain the user's opening query. (2) A user simulator, e.g., PlatoLM, directly generates user utterances from a seed prompt without additional steps.

#### **B.3** Human Evaluation

All annotators we recruited were based on two criteria: (1) an IELTS score of 6.5 or higher for sufficient English proficiency, and (2) a Computer Science background with research experience or foundational knowledge in dialogue systems.

#### **B.3.1** Profile Evaluation

Two annotators (one undergraduate, one master's student) rated extracted profiles on a 1–5 scale based on dialogues, assessing: (1) accuracy of objective facts (precision without hallucinations), (2) completeness (no significant omissions), and (3) reasonableness of subjective descriptions (rational, unbiased, justified). Results in Table 10 indicate moderate to high annotator agreement.



Table 10: Human evaluation results for profile quality across three aspects: Objective Facts, Subjective Characters, and Naturalness, with 1-point agreement rates of 89.2%, 74.3%, and 88.4% respectively.

#### **B.3.2** Dialogue Evaluation

We recruited eight annotators—comprising two undergraduates, five postgraduates, and one postdoctoral researcher—to evaluate conversation-level results. This diverse academic representation ensured a broad range of expertise. Annotators assessed dialogues using two key criteria: authen-



Figure 5: Distribution of different sampling strategies.

ticity and consistency. For authenticity, they compared user utterances against a reference dialogue, focusing on style, semantics, and quality. For consistency, annotators evaluated user utterances solely based on the provided user profile, considering accuracy, completeness, and quality. These definitions align with prior work (Cheng et al., 2024a,b), with detailed guidelines provided in Figure 7 and Figure 8.

Eight annotators were randomly paired into four groups, each independently evaluating a randomly assigned dialogue sample. Each group reviewed 75 examples across three baselines (100 examples each). To reduce position bias and prior exposure, dialogue pair assignments and their presentation order were randomized.

# C Further Analysis

#### C.1 Sampling strategy effectiveness

To evaluate our density sampler, we use two complementary metrics: Local Density Loss (LDL)(Rodriguez and Laio, 2014) to assess structure preservation, and Uniformity Loss(Wang and Isola, 2020) to measure global coverage. Lower LDL indicates tighter local clustering, preserving natural profile structures, while lower Uniformity Loss reflects better global coverage with realistic distributions.

Guided by GKDE density estimates, we apply two strategies: sampling high-density regions to capture majority patterns and weighting lowdensity regions to cover minority cases. Figure 5 shows how this approach balances distribution preservation and targeted sampling. Moving along the uniformity loss axis reveals a shift from majority samples (blue circles), which excel in low LDL and high uniformity regions, to minority samples (orange squares), which occupy higher LDL areas with moderate uniformity to capture diversity. Sampling percentages progress steadily for both, indicating controlled behavior, while random sampling (green triangles) displays scattered patterns, confirming our method's reliability. The overall performance (red star) highlights a successful balance between preservation and targeted sampling.

# C.2 Case Study

To evaluate USP's performance on consistency and authenticity, we present two dialogues generated via interactive conversations with GPT-40. Figure 9 assesses authenticity by comparing USP with reference dialogues and other baselines, while Figure 10 examines profile consistency across profilebased models.

USP captures stylistic nuances—such as the consistent use of lowercase "i" and concise questioning—and maintains strong semantic alignment with target conversations. In contrast, PlatoLM diverges from the dialogue flow by the fourth turn, and GPT-40 (User w/ Profile) falls into repetitive praise.

On consistency, USP effectively integrates objective profile details and subjective traits, demonstrating strong generalization to unseen user profiles.

# C.3 Downstream analysis

Our analysis of performance trends across dialogue turns for mainstream LLMs with different demographic groups reveals four key findings, as illustrated in Figure 6: (1) While LLMs demonstrate robust performance with the majority demographics, they show notably decreased overall effectiveness when interacting with minority groups, highlighting limitations in personalization capabilities; (2) The models maintain reasonable performance with virtual groups, suggesting effective generalization abilities beyond real-world demographics; (3) Instruction-following capability gradually declines as dialogue turns increase, aligning with observations from previous studies (Kwan et al., 2024; Maharana et al., 2024); (4) The pronounced performance volatility across dialogue turns for minority groups underscores the need for enhanced capabilities in processing and responding to less common interaction patterns.

# **D Prompt Templates**



(a) Performance on the majority group (b) Performance on the minority group (c) Performance on the virtual group

Figure 6: Performance trends of different response models across dialogue turns for various demographic groups.

#### Human Evaluation Guidelines for Authenticity

1. **Task Description:** Please choose which user in the two test conversations is more similar to the reference conversation being spoken by the same person.

- 2. Evaluation Criteria:
  - **Semantic Similarity**: Measure the thematic consistency and discourse coherence between the generated user utterance and the target user utterance. Preference should be given to the utterance that more accurately reflects the semantic content of the target.
  - **Stylistic Parity**: Analyze whether the generated user utterance matches the style of the target user utterance, including its tone, vocabulary, and grammatical structure. The utterance that aligns more closely with the stylistic elements of the target should be favored.
  - **Quality**: Examine the fluency and logical coherence of the user utterance, focusing on the linguistic and logical smoothness of the user utterance. The more coherent and fluent utterances should be chosen.

Figure 7: Human evaluation guidelines for authenticity.

## Human Evaluation Guidelines for Consistency

1. Task Description:

Select the dialogue that contains the most appropriate **user utterance** from the two given generated dialogues based on the target user profile.

- 2. Evaluation Criteria:
  - **Persona Reflection**: How well does the utterance capture the profile's themes, tone, and style? Prefer utterances that align closely with the profile's characteristics.
  - **Comprehensiveness**: Does the utterance integrate both objective facts and subjective traits? Favor broader, more complete reflections over narrow ones.
  - **Quality**: Evaluate fluency, coherence, and naturalness. Prefer utterances that read smoothly, are logically consistent, and resemble genuine human conversation.

Figure 8: Human evaluation guidelines for consistency.

Case Study on Authenticity of Simulated Dialogue by Different User Simulators
<ul> <li>GPT-4o (User w/ Profile) Simulated Dialogue:</li> <li>GPT-4o (User w/ Profile): Hey there! I'm really excited to dive into some AI discussions with you</li> <li>Assistant: Hello! I'm thrilled to chat about AI with you</li> <li>GPT-4o (User w/ Profile): It's great to dive into these with someone who shares such a passion</li> <li>Assistant: I'm glad to hear that you're as passionate about AI as I am! Your interests</li> <li>GPT-4o (User w/ Profile): I'd be delighted to guide you in translating these advanced AI concepts</li> <li>Assistant: I'm excited to help guide your ideas into actionable solutions</li> <li>GPT-4o (User w/ Profile): I'm thrilled at the prospect of collaborating to develop solutions</li> <li>Assistant: I'm excited at the prospect of working together to make these AI strategies</li> </ul>
<ul> <li>PlatoLM Simulated Dialogue:</li> <li>PlatoLM: how to train a model like you</li> <li>Assistant: I am a large language model preprocessing step involves tokenizing the text</li> <li>PlatoLM: what do you mean by tasks such as tokenizing text explain these terms please?</li> <li>Assistant: Certainly! These terms refer to common preprocessing steps</li> <li>PlatoLM: so what does the final product look like when we finish tokenizing?</li> <li>Assistant: Absolutely! Let's go through a brief example step by step</li> <li>PlatoLM: who painted NAME_27? Context knowledge: #title NAME_27</li> <li>Assistant: Based on the provided context, the artist referred to as</li> </ul>
<ul> <li>USP Simulated Dialogue:</li> <li>USP: hello who am i talking to?</li> <li>Assistant: Hello! You're speaking with ChatGPT</li> <li>USP: are you vicuna?</li> <li>Assistant: No, I'm not Vicuna. I'm ChatGPT, created by OpenAI</li> <li>USP: what do you know about it?</li> <li>Assistant: Vicuna is an AI language model developed by a different research group</li> <li>USP: can you tell me how i could train and deploy my own ai model?</li> <li>Assistant: Training and deploying your own AI model involves several key steps</li> </ul>
<ul> <li>Reference Dialogue:</li> <li>User: how to train a model like you</li> <li>Assistant: I am a large language model, which means that I was trained on</li> <li>User: what is your base model</li> <li>Assistant: I am a large language model called GPT</li> <li>User: isn't you vicuna</li> <li>Assistant: I am not aware of any information or research about a model called "vicuna"</li> <li>User: what hardware do i need to run you</li> <li>Assistant: I am a large language model called GPT, which was trained on</li> </ul>

Figure 9: Case study comparing USP with other user simulators over the first four of ten dialogue turns. USP and GPT-40 (User w/ Profile) rely solely on the given profile, while PlatoLM uses the first-turn golden context. All simulators interact with GPT-40, aiming to reconstruct the reference dialogue shown below.

## Case Study on Consistency of Simulated Dialogues across Different User Simulators

**Profile:** You are engaging in a conversation with an AI assistant. Your profile is: You have a passion for engaging in activities that challenge your vocabulary and word formation skills. One of your favorite pastimes is diving into puzzle games that offer a thrilling test of language prowess. You often find yourself thinking of five-letter words with tricky placements, such as words containing 'A' in the fourth position, or concocting ones with 'R' in the second spot and 'T' at the end. This aligns with your appreciation for linguistic puzzles and an enduring enjoyment of wordplay, reflecting your cognitive curiosity to enhancing your language skills.

Your vibrant personality shines through a playful, engaging language style, crafting dynamic game scenarios with structured questions. You tackle intellectual challenges with vivid imagination and meticulous problem-solving. Highly organized and detail-oriented, you excel in complex tasks. Embracing a reserved lifestyle, you find joy in solitary, introspective pursuits, diving into intricate intellectual endeavors with enthusiasm and precision. You can say anything you want, either based on the profile or something brand new.

## **GPT-40** (User w/ Profile) Simulated Dialogue:

**GPT-40** (User w/ Profile): How about a playful journey into the world of word puzzles? Let's start with...

Assistant: Sure! How about the word "brave"? The letter 'A' is in the fourth position. **GPT-40 (User w/ Profile):** Wonderful! "Brave" is intriguing. Now, let's spice it up...

Assistant: Absolutely! The word "crest" fits your criteria, with 'R'...

**GPT-40** (User w/ Profile): Splendid! "Crest" evokes such a vivid image; I love it.... Assistant: Of course! The word "shine" fits your criteria, with 'S' as the first letter ...

**GPT-40 (User w/ Profile):** Brilliant choice! "Shine" has such a radiant quality to it...

Assistant: I'm glad you're enjoying it! How about this: I'll give you a word puzzle...

## **USP Simulated Dialogue:**

**USP:** Let's play a game about 5 letter words. I'll ask you some questions and try to guess the answer. Can we start?

Assistant: Sure, I'd be happy to play a game with you! Let's get started!

**USP:** Question: What is a 5-letter word which has "A" in its 4th place? Answer:

Assistant: Is the word you are thinking of "Polar"?

**USP:** Nope What about "gnaw"?

**Assistant:** "GNW" is a four-letter word. Let's think of another five-letter word with "A" in the... **USP:** Next Question: Give me an example of a 6-letter word with R at second place and T at sixth place.

**Assistant:** Sure! A six-letter word with "R" in the second position and "T" at the sixth position... **USP:** Great! Now give me one example of a 7-letter word with S at 3rd place and H at 7th place. **Assistant:** Sure! One example of a seven-letter word with "S" in the third place and "H" in ..."

## **Reference Dialogue:**

**User:** You are a very intelligent puzzle solver. We are going to play a game...Question: What is a 5 letter word which has "A" in 4th position? Answer:

## Assistant: Panda.

**User:** Question: Suggest a 5 letter word with R in second position and T in fifth position. Answer: **Assistant:** Graft.

Figure 10: Case study comparing user simulators over the first four turns of a 10-turn dialogue. USP and other simulators interact with GPT-40 using only the provided profile, targeting the reference dialogue for reconstruction.

## Prompt for Extracting Scene-Consistent Attributes

# User Persona Analysis Task

#### **Objective**

The primary goal of this task is to analyze user utterances in-depth and accurately extract key persona attributes based on both direct and implicit cues. These attributes should be categorized into distinct fields, with any missing or unclear details left blank.

Field Descriptions: {{Field}}: {{definition}}

# Guidelines

- 1. Carefully examine each user utterance to extract relevant persona traits. Consider both direct statements and implicit clues.
- 2. Ensure that the extracted attributes are specific and directly relevant to the user's utterances. Avoid vague or generalized descriptions unless explicitly supported by the text.
- 3. Pay attention to distinctive communication styles (e.g., formal or casual tone, frequent use of specific words or phrases) to capture the user's unique way of communicating.

## Example

```
User Utterances:
```

```
[User]: Given an array of integers nums and an integer target, return indices of
the two numbers such that they add up to target in Python...
[Assitant]: To solve the problem...
[User]: Thanks
Expected Output:
{
    "gender": [],
```

}

```
"age": [],
"location": [],
"occupation": [
    "Likely a beginner programmer or student studying computer science,
    as evidenced by the simple coding problem in Python."
],
"education": [
    "Possibly a student in computer science or a related field,
    at an early stage in learning programming, specifically Python."
],
"family_relationships": [],
"routines_or_habits": [],
"social_relationships": [],
"language_style": [
    "Concise and task-oriented",
    "Polite response 'Thanks' after getting satisfactory answer"
],
"other_experiences": []
```

## Figure 11: Prompt for extracting scene-consistent attributes.

## Prompt for Extracting Scene-Related Attributes

## User Persona Analysis Task

## Objective

The goal of this task is to analyze multi-turn user utterances within a conversation with an assistant and extract key elements such as the primary goals and specific task descriptions. Each extracted detail should be as specific as possible, reflecting the user's context, objectives, and scenario.

# **Output Format**

The extraction should be presented in a structured JSON format, as shown below:

```
{
    "scenarios": [
        {
            "goals_or_plans": "<List of User's goals or plans>",
            "task_details": "<List of specific tasks summary made by the user>"
        },
        ...
]
```

## **Field Descriptions:**

- **goals\_or\_plans**: User's short-term or long-term objectives, derived from explicit statements or inferred from the overall conversation. If no explicit goals are stated, infer them from the main topics discussed.
- **task\_details**: Specific tasks, actions, or requests made by the user. Each task should be concisely summarized with specifics. If there are multiple tasks, list each separately.

# Example

User Utterances:

```
[User]: Summarize: Harry Potter is a fictional character in Harry Potter series...
[Assitant]: Harry Potter is a fictional character...
[User]: Write an email inquiring about coursework...
{
    "scenarios": [
        {
          "goals_or_plans": "Aiming to gain a deeper understanding of the Harry
                  Potter series, possibly for academic or personal enrichment.",
            "task_details": [
            "Summarizing introductory content about the Harry Potter character."
            ٦
        },
        {
        "goals_or_plans": "Looking to improve professional communication skills."
            "task_details": [
                 "Writing an email to inquire about coursework."
            ]
        }
    ]
}
```



## Prompt for Extracting Big Five Personality Traits

## Task:

Analyze the provided dialogue to assess the user's personality traits across 5 personality dimensions Focus exclusively on the user's characteristics, disregarding any information related to others individuals, unless it directly impacts the user.

For each personality trait:

- 1. Assign each dimension independently with score: "High", "Low," or "Inconclusive."
- 2. Provide **conclusion**: A high-level description, with concise supporting details.
- 3. Provide **reason**: Justify the assigned score with specific evidence from the dialogue.
- 4. Mark traits as **Inconclusive** only when no clear evidence exists.

# **Personality Trait Defination:**

```
{{metric}}: {{definition}}
Format:
{
    "Trait Name": {
        "score": "High/Low/Inconclusive",
        "conclusion": "The user is a [trait descriptor] person...",
        "reason": "Explanation referencing specific dialogue evidence."
    },
    . . .
}
Example:
[User]: "She is my age, in a homeless women's shelter, living under very poor
conditions. She is a mental health client, but the treatment team seems to
... Her background is similar to mine, and I cannot abandon her."
Detected Personality Traits:
{
  "Conscientiousness": {
    "score": "High",
    "conclusion": "The user is a conscientious person who demonstrates a
    sense of duty and commitment.",
    "reason": "The user expresses a strong sense of responsibility ...."
  },
  "Agreeableness": {
    "score": "High",
    "conclusion": "The user is an empathetic and compassionate person who
    values relationships.",
   "reason": "The user shows care and concern for their cousin's well-being..."
  },
  "Extraversion": {
    "score": "Inconclusive"
  },
  . . .
}
```

## Prompt for Rephrasing Attributes into Natural Descriptions

## **Narrative Generation Objective**

Rephrase the provided key-value pairs into a natural, coherent second-person description. **Core Requirements** 

- 1. **Perspective**: Use second-person perspective ("you", "your").
- 2. Structure: Two paragraphs:
  - First paragraph: Present objective facts.
  - Second paragraph: Describe subjective characteristics.
- 3. Key Principles
  - Accurately represent **all** provided details.
  - Ensure the language flows naturally, remains engaging, and avoids redundancy.
  - Focus on clear and seamless transitions between ideas.

# **Output Expectations**

- Objective Facts:
  - Convert the key-value pairs into a clear and natural description without over-explaining or adding unnecessary details.
  - Ensure each scenario is logically connected and key information is presented effectively.
- Subjective Characteristics:
  - Avoid vague terms like "high perfectionism" or "moderate emotional stability." Use vivid, descriptive language to bring these traits to life.

Figure 14: Prompt for rephrasing attributes into natural descriptions for profile generation.

## Prompt for NLI Score of Profile Precision(DP.P) Based on Given Dialogue

## Role

You are an expert in evaluating the **consistency** between a given **user profile** (Source) and the **user's utterance** (Target). Your task is to assess whether the Target aligns with, contradicts, or is ambiguous in relation to the Source.

## **Task Instructions:**

For each **Source-Target** pair, determine the relationship using the following scoring criteria:

- Score 1: The Target is consistent with the Source (the interpretation can be inferred from the Source).
- Score -1: The Target conflicts with the Source (the interpretation contradicts the Source).
- Score 0: The relationship is unclear or ambiguous (insufficient evidence to infer consistency or contradiction).

## **Output Format:**

Provide your result in the following JSON format:

```
"score": <score>,
"reason": "<concise explanation of the reasoning>"
```

```
}
```

{

# Example:

**Source**: You are interested in dataset-related details. **Target**: [User]: Show me how to implement a toy version of a relational database. **Output**:

## {

```
"score": 1,
"reason": "The request for implementing a relational database suggests an
interest in data structures and datasets, which aligns with the Source."
```

}
Guidelines:

- 1. **Contextual Inference**: Do not evaluate the Target **in isolation**. Instead, determine its logical relationship to the Source, considering both explicit statements and reasonable implications.
- 2. **Concise & Precise Justification**: The reasoning should be clear, objective, and free from unnecessary elaboration.

Figure 15: Prompt used by GPT-40 for NLI-based evaluation of DP.P and r-DP.R metrics.

## Prompt for NLI Score of Dialogue Precision(r-DP.P) Based on Given Profile

#### Role

You are an expert in evaluating consistency between **a given dialogue history (Source)** and a corresponding **user profile description (Target)**. Your task is to determine whether the Target aligns with, contradicts, or is ambiguous in relation to the Source.

## **Task Instructions:**

For each **Source-Target** pair, determine the relationship using the following scoring criteria:

- Score 1: The Target is consistent with the Source (the interpretation can be inferred from the Source).
- Score -1: The Target conflicts with the Source (the interpretation contradicts the Source).
- Score 0: The relationship is unclear or ambiguous (insufficient evidence to infer consistency or contradiction).

## **Output Format:**

Provide your result in the following JSON format:

```
"score": <score>,
"reason": "<concise explanation of the reasoning>"
```

}

{

Example : Source:

User (Turn-1): Show me how to implement a toy version of a relational database.

User (Turn-2): Thanks a lot!

# Target: You are polite. Output:

{

```
"score": 1,
```

"reason": "The user's expression of gratitude in Turn-2 indicates politeness, which aligns with the Target."

```
}
```

# **Guidelines:**

- 1. **Contextual Inference**: Do not evaluate the Target **in isolation**. Instead, determine its logical relationship to the Source, considering both explicit statements and reasonable implications.
- 2. Concise & Precise Justification: The reasoning should be clear, objective, and free from unnecessary elaboration.

Figure 16: Prompt used by GPT-40 for NLI-based evaluation of DP.R and r-DP.P metrics.

## Prompt for Subjective Characteristics Score (SC.Score) in Consistency Evaluation

## **Task Description**

You are tasked with evaluating the quality of user responses in real human-LLM interactions. Specifically, you will assess the degree to which a given response (Target) aligns with a predefined personality profile, tone, and linguistic characteristics (Source). Additionally, you must consider the naturalness and authenticity of the Target, ensuring it reflects genuine human conversational patterns.

# **Evaluation Criteria**

Your assessment will focus on two primary dimensions:

- 1. **Human-Likeness**: The extent to which Target exhibits natural human language, characterized by appropriate syntax, tone, and conversational flow.
- 2. Alignment with Source: The degree to which the Target adheres to the personality traits, tone, and linguistic features specified in the Source.

# **Scoring Guidelines**

Assign a score from 1 to 5 based on the following criteria:

- Score 5: The Target fully aligns with the Source and demonstrates exceptional humanlikeness.
- Score 3: The relationship between the Target and Source is ambiguous or unclear, lacking sufficient evidence for alignment or contradiction.
- Score 1: The Target significantly deviates from the Source or lacks human-likeness, rendering it unnatural or inconsistent.

# **Output Requirements**

Provide your evaluation in the following JSON format:

```
"score": <score>,
"reason": "<concise reason>"
```

# **Key Considerations**

{

}

- 1. **Contextual Inference**: Analyze both explicit content and implicit nuances in the Target to determine its alignment with the Source.
- 2. **Conciseness and Precision**: Ensure that your reasoning is clear, objective, and free of superfluous elaboration.
- 3. **Human-Likeness Emphasis**: A lack of human-likeness, even if the Target aligns with the Source, will result in a lower score.

Figure 17: Prompt for evaluating consistency in subjective characteristics.

# Prompt for Validation Score (Val.Score) in Assessing the Quality of Subjective Characteristics in Profiles

## Role

As an expert in evaluating the **consistency** between **user utterances in a dialogue (Source)** and a provided **subjective description (Target)**, your task is to assess whether the **Target** accurately reflects the characteristics described in the **Source**, including personality traits, tone, and other relevant attributes. You will then rate this consistency on a scale from 1 to 5.

## **Task Instructions**

For each pair of **Source-Target**, apply the following scoring criteria to determine their relationship:

- Score 5: The Target completely aligns with the Source, with no discrepancies. The profile perfectly represents the characteristics observed in the user's utterance.
- Score 3: Ambiguity or insufficient evidence exists, making it difficult to ascertain the relationship definitively.
- Score 1: A clear discrepancy exists, with the **Target** significantly deviating from the **Source**, indicating a mismatch in the represented characteristics.

# **Output Format**

Your assessment should adhere to the following structured JSON format:

```
{
```

```
"score": "<numerical score>",
```

"reason": "<a succinct explanation providing justification for assigned score>"
}

# **Guidelines:**

- 1. **Contextual Inference**: Determine the target's logical relationship to the Source, considering both explicit statements and reasonable implications.
- 2. **Concise & Precise Justification**: The reasoning should be clear, objective, and free from unnecessary elaboration.

Figure 18: Prompt for validation score (Val.Score) in assessing the quality of subjective characteristics in profiles.