SIMUSER: WHEN LANGUAGE MODELS PRETEND TO BE BELIEVABLE USERS IN RECOMMENDER SYSTEMS

Anonymous authors

Paper under double-blind review

Abstract

Recommender systems play a central role in numerous real-life applications, yet evaluating their performance remains a significant challenge due to the gap between offline metrics and online behaviors. We introduce SimUSER, an agent framework that serves as believable and cost-effective human proxies for the evaluation of recommender systems. Leveraging the inductive bias of foundation models, SimUSER emulates synthetic users by first identifying self-consistent personas from historical data, enriching user profiles with unique backgrounds and personalities. Then, central to this evaluation are users equipped with persona, memory, perception, and brain modules, engaging in interactions with the recommender system. Specifically, the memory module consists of an episodic memory to log interactions and preferences, and a knowledge-graph memory that captures relationships between users and items. The perception module enables visualdriven reasoning, while the brain module translates retrieved information into actionable plans. We demonstrate through ablation studies that the components of our agent architecture contribute to the believability of user behavior. Across a set of recommendation domains, SimUSER exhibits closer alignment with genuine humans than prior state-of-the-art, both at micro and macro levels. Additionally, we conduct insightful experiments to explore the effects of thumbnails on click rates, the exposure effect, and the impact of reviews on user engagement. The source code is released at https://github.com/SimUSER-paper/SimUSER.

032

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

Recommender systems (RS) have become an indispensable component of our day-to-day lives, of-033 fering personalized user experience and improving satisfaction Li et al. (2024b). By examining 034 granular preferences, historical data, and contextual factors, RS can deliver recommendations that 035 are tailored to meet individual tastes. Despite their widespread adoption, a key challenge hindering the advancement of the field is evaluation Yoon et al. (2024). The difficulty arises from the dis-037 crepancy between offline metrics (non-interactive), which are typically used during development, and real-life user behaviors, which these systems encounter post-deployment Zhang et al. (2019). This results in models that perform well in controlled environments but fail to meet expectations in 040 practical use cases. Such a limitation is further exacerbated by the inherent shortcomings of offline 041 evaluation, notably the inability to measure business values such as user engagement and satisfac-042 tion Jannach & Jugovac (2019). On the other hand, deploying and testing RS in real-world settings 043 is both costly and labor-intensive, underscoring the imperative need for reliable and affordable (interactive) evaluation methods. 044

Recent breakthroughs in Large Language Models (LLMs) have shown promise in human behavior modeling by enabling the creation of autonomous agents. Generative agents have demonstrated capabilities in lifelong learning through automatic skill discovery Wang et al. (2023a) or setting consistent goals during exploration Du et al. (2023). LLMs have also been applied to various social simulations, such as hospital Li et al. (2024a) and city Park et al. (2023). In a different spirit, S³
Gao et al. (2023) simulates the dynamic evolution of opinions in a social network. In the realm of recommendation systems, RecMind Wang et al. (2023e) explores the concept of autonomous recommender agents equipped with self-inspiring planning and external tool utilization. Recently, InteRecAgent Huang et al. (2023) has extended this idea by proposing memory components, dynamic demonstration-augmented task planning, and reflection. Nevertheless, their primary focus 054 remains on identifying candidate items that align with user preferences. In addition, these agents 055 are overly simplified representations of human users, limited to predicting item ratings based on 056 viewing history. Recently, RecAgent Wang et al. (2023c) has attempted to introduce more diverse 057 user behaviors, taking into account external social relationships. Another work, Agent4Rec Hou 058 et al. (2024), delves into generating faithful user-RS interactions via agent-based simulations, where agents are equipped with a memory module. However, a common characteristic of existing studies is their insulated nature — they primarily rely on knowledge embedded within the model's 060 weights, neglecting the potential benefits of integrating external knowledge and user-item relation-061 ships. Furthermore, prior approaches often disregard user personas (e.g., personality traits) and fail 062 to incorporate visual signals, despite images significantly shaping user experience and emotion. 063

To enable synthetic users, this paper describes an agent architecture built upon large language mod-064 els. Our methodology consists of two phases: (1) self-consistent persona matching and (2) recom-065 mender system evaluation. In Phase 1, we leverage the semantic awareness of LLMs to extract and 066 identify consistent personas from historical data, encompassing unique backgrounds, personalities, 067 and characteristics. In Phase 2, we impersonate these identified personas to simulate believable hu-068 man interactions. This involves a retrieval-augmented framework where the agent interacts with the 069 recommender system based on its persona, memory, perception, and brain modules. The memory module comprises an episodic memory and a knowledge-graph memory. The former is the episodic 071 memory that records, in natural language, a comprehensive list of the agent's experiences. The lat-072 ter captures relationships between users and items, accounting for the influence of other users and 073 prior beliefs about items. Unlike existing studies that solely rely on text, our perception module 074 incorporates visual cues into the agent's reasoning process. Finally, the brain module is respon-075 sible for translating retrieved information and the current simulator observation into action plans such as watch, click, or exit. Following action selection, the user engages in self-reflection to 076 synthesize memories into higher-level inferences and draw conclusions. 077

To our knowledge, this study is one of the pioneering works in developing a general LLM-based agent framework for systematic evaluation of recommender systems, integrating persona extraction, visual-driven reasoning, and domain-specific prior knowledge to simulate realistic user interactions.

081 082

083

084

2 RELATED WORK

085 LLMs as Human Proxies. The use of LLMs to replace costly human trials is becoming increasingly popular. For instance, Owoicho et al. (2023) and Wang et al. (2024d) examine LLMs as user 087 simulators for conversational search, while Hämäläinen et al. (2023) seeks to generate synthetic user 880 experience data in human-computer interaction (HCI). Significant progress have also been achieved in the area of collaborative cooperation. For instance, CAMEL Li et al. (2023a) describes a frame-089 work for agent cooperation that features a commander for planning and executors for task imple-090 mentation. Qian et al. Qian et al. (2023) introduces a virtual software company where agents, each 091 assigned roles such as engineer, work together to complete software development projects. Besides, 092 Aher et al. (2023) explores LLMs' ability to replicate human behavior for social science experi-093 ments. Other researchers have created simulation environments where LLM agents interact with 094 one another, producing realistic daily activities Park et al. (2023); Gao et al. (2023) or generating 095 mobility trajectories Wang et al. (2024b). Nonetheless, their potential for assessing recommender 096 systems and addressing inherent issues in recommendation remains largely unexplored.

Simulating Users in Recommendation. Conversational RS initially tackled the recommendation 098 problem using bandit models, emphasizing the quick update of traditional systems through item selection and binary feedback from synthetic users Christakopoulou et al. (2016). Taking this further, 100 Zhao et al. (2023) created a simulation platform where users not only chat about recommendations 101 but also navigate websites, search for items, and share opinions on social media. Recent techniques 102 have added more natural language flexibility, but user responses are usually limited to binary or 103 multiple-choice formats Lei et al. (2020). Evaluations typically involve predefined target items for 104 each user, with success measured by the number of turns needed to identify these targets Guo et al. 105 (2018); Sun & Zhang (2018). Another technique employs agenda-based simulations, utilizing state diagrams to guide actions, and deeming a recommendation successful when the conversation reaches 106 a "complete" state Zhang & Balog (2020); Zhang et al. (2022). In spite of this, these simulations 107 often rely on fixed rules and scripted dialogues, lacking the variability seen in human interactions.

108 LLMs in Recommender Systems. To address the above-mentioned limitations, generative sim-109 ulators using Large Language Models (LLMs) have been developed, offering more realistic and 110 nuanced conversational abilities Zhang et al. (2024b); Zhao et al. (2023). A few studies have also 111 explored the application of LLMs as recommender systems Hou et al. (2024); Li et al. (2023b); 112 Kang et al. (2023); Fan et al. (2023); Yang et al. (2023a). These investigations explore LLMs as recommendation engines, rather than as entities that perceive recommendations, thus providing a 113 perspective complementary to our research Wang et al. (2024c); Zhang et al. (2024a). Some authors 114 Wang et al. (2023d) utilize ChatGPT to prompt users towards specific items, where users would 115 provide 'hints' about these items. LLM4Rec Wang et al. (2024a) alleviates the beam search decod-116 ing by using a straight item projection head for ranking scores generation, providing speedup when 117 predicting the next item. Rec-GPT4V investigates the use of vision-language models for multimodal 118 recommendation Liu et al. (2024). In order to bridge the gap between reasoning capabilities and ex-119 ternal knowledge, KAR Xi et al. (2023) incorporates open-world knowledge, which includes factual 120 knowledge on items and user preferences. An emerging line of work seeks to simulate user behav-121 iors in the context of RS. For instance, RecAgent Wang et al. (2023b) presents a new simulation 122 framework aimed at enhancing the evaluation of recommendation algorithms through agent-based 123 simulation. In a similar fashion RecMind Wang et al. (2023e) proposes self-inspiring agents for recommendation. However, their simulated users are limited to basic actions like rating items, lack-124 ing the ability to engage in more complex interactions, such as navigating the interface. Notably, a 125 recent approach Yoon et al. (2024) examines the effectiveness of LLMs as generative users, specif-126 ically for conversational recommendation scenarios. A closely related work to ours is Agent4Rec 127 Zhang et al. (2023) that delves into the generative capabilities of LLMs for modeling user interac-128 tions. SimUSER differs significantly from these studies in several key aspects. We simulate realistic 129 humans using detailed personas that are inferred from their interaction history. We also incorporate 130 a perception module to integrate visual information in decision-making, mapping thumbnail items 131 to token embeddings of the language model. Furthermore, the present study explores the potential 132 for graph-based retrieval and features a richer spectrum of actions such as "going to previous page" 133 or "clicking", to better replicate influential factors and user experience.

134 135

3 Methodology

136 137

138 Simulated USERs provides a framework for 139 systematically assessing recommender systems 140 by engaging in interactions and providing feed-141 back. As illustrated in Figure 1, we present 142 a two-phase method. Phase 1 of SimUSER matches historical data with a set of personas to 143 enable nuanced and realistic interactions with 144 the recommender system. Phase 2 utilizes the 145 identified personas, historical data, and novel 146 reasoning mechanisms to generate synthetic 147 users with human-like behavior. At the cen-148 ter of our architecture is the memory module. 149 From the memory module, interaction records 150 and prior assumptions are retrieved to condi-151 tion the agent's actions and react appropriately



Figure 1: Illustration of SimUSER framework.

to the simulation. Additionally, the perception module enables agents to transform visual cues into a natural language form when interacting with the recommender system.

Problem Formulation. Given a user $u \in \mathcal{U}$ and an item $i \in \mathcal{I}$, the aggregated rating of the item is denoted by $R_i = \frac{1}{\sum_{u \in \mathcal{U}} y_{ui}} \sum_{u \in \mathcal{U}} y_{ui} \cdot r_{ui}$ where $y_{ui} = 0$ indicates that the user u has not rated the item i and inversely $y_{ui} = 1$ indicates that the user has rated the item with $r_{ui} \in \{1, 2, 3, 4, 5\}$. We also introduce $g_i \in G$ as the genre/category of the item. In this study, we seek to discover y_{ui} and r_{ui} for an unseen recommended item i.

160 Generative Large Language Models. LLMs are trained to predict the most probable next to-161 ken t_k given the sequence of previous tokens $t_1 ldots t_{k-1}$ by maximizing the likelihood function $p_{LLM}(t_k|t_1, \dots, t_{k-1})$. In this work, we use pre-trained LLMs without further finetuning them. ¹⁶² Depending on the task, we generate one or more tokens given a task-specific context $c^{(p)}$ that describes the task to the language model and prompts it for an answer. Thus, we obtain generated tokens **t** by sampling from:

165 166 167

168 169

170

$$p_{LLM}(t|c^{(p)}) = \prod_{k=1}^{K} p_{LLM}(t_k|c_1^{(p)}, \dots, c_n^{(p)}, t_1, \dots, t_{k-1})$$
(1)

3.1 PERSONA MATCHING VIA SELF-CONSISTENCY CHECK

This phase involves assessing the most plausible *persona* based on historical data. A persona *p* encompasses a set of features that characterize the user: **age**, **personality**, and **occupation**. Personality traits are defined by the Big Five personality facets: *Openness*, *Conscientiousness*, *Extraversion*, *Agreeableness*, and *Neuroticism*, each measured on a scale from 1 to 3. We postulate that these factors are critical for modeling personalized preferences. Given the difficulty of obtaining such granular features in real-world settings, our methodology seeks to systematically infer personas from the user's interaction history.

178PersonaExtraction.Assuming a user u and a set of interactions179 $\{(i_0, r_{ui_0}), (i_1, r_{ui_1}), \dots, (i_n, r_{ui_n})\}$, we first query the LLM to produce a short summary180 s_u of the user's preferences. To do so, we randomly select 50 items from the user's viewing history.181Items rated 4 or above are categorized as *liked*, while those rated below 3 are deemed *disliked*. This182summary describes its unique tastes and rating patterns.

Following this, both the summary s_u and historical data are combined in the prompt, instructing the LLM agent to generate a persona that matches the interaction history for this user. To enhance the diversity in the generated personas, the LLM is provided a list of possible ages, personalities, and occupations. For each user, a set of m (m = 5) candidate personas is generated, denoted as \mathcal{P} .

Self-Consistent Persona Evaluation. We then assess the consistency of the candidate personas \mathcal{P} to identify the most plausible one. Specifically, a self-consistency scoring mechanism measures the alignment of candidate personas with historical data. We define a scoring function s(p, u) for each candidate persona $p \in \mathcal{P}$, where p is evaluated against two distinct sets of user-item interactions. For the targeted user u, we sample j subsets of ϱ interactions from its history. These are compared with ϱ sampled interactions from other users \bar{u} , denoted as $I_{\bar{u}}$:

194 195

202

204

 $s(p,u) = \sum_{\iota \sim I_u} \hat{r}(\iota, p) - \sum_{\bar{\iota} \in I_{\bar{u}}} \hat{r}(\bar{\iota}, p)$ (2)

where $\hat{r}(\iota, p)$ and $\hat{r}(\bar{\iota}, p)$ are obtained by querying the LLM to rate the two interaction subsets ι and i. Ideally, the LLM agent should assign a higher $\hat{r}(\iota, p)$ for interactions from the targeted user and a lower $\hat{r}(\bar{\iota}, p)$ for samples from other users. This approach ensures self-consistency in the persona: the persona derived from the rating history should be consistent with the data during the evaluation. The candidate persona p with the highest score s is assigned to the user. We provide the prompts for Phase 1 in Appendix B.

203 3.2 ENGAGING IN INTERACTIONS WITH RECOMMENDER SYSTEMS

In Phase 2, given a user u and its discovered persona p, we present a novel cognitive architecture built upon LLMs to simulate faithful human proxies. This model accounts for various factors, including the user's persona, memory of past interactions, habits, and unique tastes/preferences. It comprises four major modules: **persona**, **perception**, **memory**, and **action**.

209 3.2.1 PERSONA MODULE

In the recommendation domain, the user's persona plays a central role in aligning the agent's behavior with genuine human actions Li et al. (2024b). To lay a reliable foundation for the generative agent's subsequent interactions and evaluations, benchmark datasets (e.g., MovieLens-1M Harper & Konstan (2015), Steam Kang & McAuley (2018), AmazonBook McAuley et al. (2015)) are used for initialization of the persona module. In detail, each agent's profile includes the matched persona p along with attributes extracted from its historical data: $p \cup {\text{pickiness}, habits, unique tastes}.$ Since LLMs are biased towards positive sentiment, unless prompted to behave as picky users Yoon et al. (2024), each agent is assigned a *pickiness* level sampled in {*not picky*, *moderately picky*, *extremely picky*} based on the user's average rating.

Habits, as defined in Zhang et al. (2023), account for user tendencies in interactions with recommender systems: engagement, conformity, and variety.

- Engagement quantifies the frequency and breadth of a user's interactions with recommended items, distinguishing between users who extensively watch and rate many of items and those who confine themselves to a minimal set. The engagement trait for user u can be mathematically expressed as: $T_{act}^u = \sum_{i \in \mathcal{I}} y_{ui}$.
 - Conformity measures how closely a user's ratings align with average item ratings, drawing a distinction between users with unique perspectives and those whose opinions closely mirror popular sentiments. For user *u*, the conformity trait is defined as: $T_{conf}^{u} = \frac{1}{\sum_{i \in \mathcal{I}} y_{ui}} \sum_{i \in \mathcal{I}} y_{ui} \cdot |r_{ui} R_i|^2$.
 - Variety reflects the user's proclivity toward a diverse range of item genres or their inclination toward specific genres. The variety trait for user u is formulated as: $T_{div}^{u} = |U_{i \in \{y_{ui}=1\}}g_i|$. To encode users' unique tastes in natural language, we utilize the summary s_u obtained in Phase 1, which describes long-term preferences.

Unique tastes of the user, derived from their viewing history, are encapsulated in the summary s_u generated in Phase 1. This summary offers a detailed overview of their preferences, including both favored and disliked genres, and watching habits.

238 239 3.2.2 PERCEPTION MODULE

222

223

224

225 226

227

228

229 230

231

232

233

234

240 A primary factor in decision-making is visual stimuli due to their significant influence on curiosity 241 and emotion Liu et al. (2024). For instance, when scrolling through a movie recommendation plat-242 form, human decisions are heavily driven by the thumbnails of items, which can trigger emotional 243 responses and provide quick visual summaries of the content Koh & Cui (2022). To graft these vi-244 sual elements in an cost-efficient manner, we propose augmenting action prompts (see Sec 3.3) with 245 image-derived captions. The caption $i_{caption}$ of an item i is generated by querying GPT-40 to extract insights that specifically capture emotional tones, visual details, and unique selling points from the 246 item's thumbnail. The captions $i_{caption} = LLM(P_{caption}, i)$ are generated only once for each item 247 before the interactions begin. For the sake of simplicity, all agents share the same captions. In future 248 work with richer user interface, we anticipate using the raw images as input for a multimodal LLM 249 or conditioning the captioning prompt $P_{caption}$ with the user's persona. 250

251 3.2.3 MEMORY MODULE

It is of vital importance for an agent to maintain a memory of the knowledge and experience it has of
the world and others. We propose to use an episodic memory storing the interactions with the RS and
a knowledge-graph memory that leverages graph-structured data to capture relationships, enabling
the model to access knowledge beyond what is inherently present in the LLM's parameters.

Episodic Memory retains a comprehensive record of the user's interactions with the recommendation system, including ratings, liked and disliked items, as well as the underlying reasons. For instance, an entry might be "You rated a movie called 'Star Wars: Episode VI' 4 out of 5". The memory is initially populated with the user's viewing and rating history. Each time SimUSER executes a new action or rate an item, the corresponding interaction is added to the episodic memory. At inference time, the agent conditions its decisions on information retrieved from this memory.

For humans, recall is the psychological process of accessing memories from the past Atkinson & Shiffrin (1968). In our system architecture, it involves retrieving "documents" from the episodic memory. We implement a self-ask strategy by prompting the LLM to raise follow-up questions regarding the query. These questions, along with the initial query, then serve as queries for vector similarity search. Given a query q, the top- k_1 documents d that have the highest similarity scores s(q, d) when compared with the query are retrieved in this step:

$$s(q,d) = \cos(\mathbf{E}(q), \mathbf{E}(d)) = \frac{\mathbf{e}_q \cdot \mathbf{e}_d}{||\mathbf{e}_q||||\mathbf{e}_d||}$$
(3)

where **E** is an embedding function.

* Knowledge-Graph Memory User behaviors in real-life recommender systems are shaped by 272 both internal and external factors Zhao et al. (2014). Examples of internal factors include user 273 personality, habits, and age, which are captured in the persona module. For external factors, we 274 consider the influence of other users and prior beliefs about items. To simulate external influences, 275 SimUSER leverages a KG memory that is used to retrieve items sharing similar relationships within 276 the graph and their intrinsic characteristics. The retrieval process emulates the way we form beliefs, 277 influenced by recommendations from friends, family, and prior knowledge, which is crucial for 278 reducing hallucinations and contextualizing user preferences. 279

280 Memory Initialization The knowledge graph memory is initially populated using real-world 281 datasets such as MovieLens, AmazonBook, and Steam. It is structured as a Knowledge Graph 282 (KG) \mathcal{G} where nodes denote entities, and edges correspond to relations between entities. An edge in 283 the KG represents a fact stored in the form of (*subject, predicate, object*). Formally, the knowledge 284 graph $\mathcal{G} = \{(h, r, t) | h, t \in \mathcal{V}, r \in \mathcal{E}\}$, in which each triple (*h,r,t*) indicates that a relation *r* exists 285 from head entity *h* to tail entity *t*. \mathcal{V} is a set of entities and \mathcal{E} represents relationships between them. 286 For instance, nodes \mathcal{V} may represent entities such as *users* and *items*, while edges \mathcal{E} may depict the 287 relationships between these entities such as *such* as *liked*.

Memory Growth The memory grows with each interaction, capturing the evolving nature of user preferences and behaviors. Formally, given an interaction i_t at time t, such as liking, disliking, or rating an item, the memory update can be represented as:

$$\mathcal{G}_{t+1} = \mathcal{G}_t \cup \{ (v_i, e_{ij}, v_j) | (v_i, e_{ij}, v_j) \in \mathcal{V} \times \mathcal{E} \times \mathcal{V} \}$$

$$\tag{4}$$

where v_i and v_j are entities (nodes) involved in the interaction i_t (e.g., user and item) and e_{ij} is the relationship (edge) created or updated by the interaction. For simplicity, this scheme assumes that agents only perceive their own interactions.

Graph-Aware Dynamic Item Retrieval For a user u, the retrieval function takes a query item x as input and returns a set of similar items along with their associated metadata (e.g., *ratings*), leveraging the KG structure and semantic similarity. A variety of algorithms may be used to measure the similarity between two items. In our pipeline, we extend Pathsim Sun et al. (2011) on account of its flexibility to capture both user-item and item-item relationships.

A relationship path $p_{x \rightarrow y}$ represents a composite relationship between entities x and y in the form of $x \xrightarrow{\mathcal{E}_1} z \xrightarrow{\mathcal{E}_2} \dots \xrightarrow{\mathcal{E}_l} y$, where \mathcal{E}_1 denotes the edge between entity x and z. For example, in the MovieLens network, the co-actor relation can be described using the length-2 relationship path $x \xrightarrow{acts-in} z \xrightarrow{actor} y$. In order to retrieve relevant items based on the query x, SimUSER estimates the item-item similarity as follows:

292

$$s_{x,y} = \frac{2 \times |\{p_{x \to y} : p_{x \to y} \in \mathcal{P}\}|}{|p_{x \to x} : p_{x \to x} \in \mathcal{P}| + |p_{y \to y} : p_{y \to y} \in \mathcal{P}|}$$
(5)

310 where \mathcal{P} denotes the set of paths between the query item x and candidate item y, and $p_{x \rightarrow y}$ is a path instance. The score $s_{x,y}$ is determined by two factors: (1) the connectivity level, which is the count 311 of paths that connect \vec{x} and y through $\hat{\mathcal{P}}$; and (2) the balance of visibility, defined by the number 312 of times these paths are traversed between the two entities. In addition to item-item similarity $s_{x,y}$, 313 we compute user-item similarity $s_{u,y}$ for the target user u and the candidate item y, using the same 314 path-based approach, which is further summed up to $s_{x,y} = \alpha \cdot s_{x,y} + (1 - \alpha) \cdot s_{u,y}$. Thus, the 315 retrieval output is influenced by the communities surrounding both the current user and the target 316 item, along with past interactions of u. 317

To account for the semantic similarity between items, the weighted cosine similarity $\frac{\mathbf{e}_x \cdot \mathbf{e}_y}{\|\mathbf{e}_x\|\|\mathbf{e}_y\|}$ is summed up with $s_{(x,y)}$, where \mathbf{e}_x is the node embedding generated with OpenAI's model *textembedding-3-small*. We found the semantic score helpful in capturing hidden variables that are not represented in the KG. Finally, the top- k_2 item candidates and their attributes are returned and serve to condition the brain module. Note that the rating of a retrieved item is substituted with the user's own rating if available; otherwise, it defaults to the average rating from other users stored in \mathcal{G} .

324 3.3 BRAIN MODULE

We endow each agent with a decision-making module that derives the agents' subsequent actions based on the retrieved information, considering their persona, tastes, emotions, and fatigue. It involves 1) **deciding** on items of interest, 2) **rating** them and **providing** feelings, 3) **selecting** the next action, and 4) **engaging** in post-interaction reflection. To replicate human-like sequential reasoning, we employ Chain-of-Thought prompting, repeatedly performing the four steps.

331 Watching Items: In our simulation, agents browse items page by page. According to their prefer-332 ences, persona, and watching history, they assess each item on the page and may decide to "watch" 333 certain items — agents select items that interest them. To address the inherent bias towards positive 334 sentiment often observed in LLMs Yoon et al. (2024), SimUSER incorporates a pickiness modifier into the prompt: You are {pickiness} about {item_type}. When available, we enrich 335 item descriptions with their thumbnail captions to enable multimodal reasoning. Finally, for each 336 item i, the top- k_2 similar items and their ratings are retrieved from the knowledge graph mem-337 ory, helping to mitigate LLM hallucinations — reducing the impact of any singular bias towards or 338 against an item, while enriching context information. 339

Providing Feelings and Rating Items This action is triggered once the user has found the items of interest. Intuitively, a real user may produce much feelings after watching an item, which will be stored in their memory and influence their future cognition and behaviors. Along with the item rating $\in \{1, 2, 3, 4, 5\}$, we query the user's feelings about the watched items and leverage such information to update the memory module. That is, newly liked and disliked items are fed back into the memory module to influence the agent's future behavior. In order to enhance consistency, users are asked to provide a reason for their feedback.

347 Emotion-driven Action Selection Once the recommender system returns the search or recommendation results, including n items each time, the agent discerns the current recommendation page and 348 retrieves its interaction history. Drawing upon these insights combined with its persona, the agent 349 decides whether to [EXIT] the system, go to [NEXT] page, return to a [PREVIOUS] page, or 350 [CLICK] on an item to access more details. To do so, the agent sequentially: 1) estimates its sat-351 isfaction level with preceding recommendations, 2) generates its current fatigue level Zhang et al. 352 (2023), 3) infers its current emotion, such as EXCITED, and 4) selects the most suitable action. 353 Satisfaction level, fatigue, and emotion are dynamic elements that the agent employs to adapt its 354 actionable plan with the recommender system. If the agent decides to click on an item, the item 355 is displayed with an extended description that replaces the short {*item_description*} as outlined 356 in Appendix B, which is then used to determine whether it wishes to engage further with the item. 357 Finally, if [EXIT] is selected, a satisfaction interview is conducted to gather granular opinions and ratings on the presented recommendations. 358

359 Post-interaction Reflection LLM-Empowered agents, when equipped with only raw episodic mem-360 ory, struggle to generalize or make inferences Park et al. (2023). Post-interaction reflection is a 361 mechanism designed to let agent learn from interactions and improve future alignment with their 362 persona. After the agent engages in assessing and rating the items on the page, we collect inter-363 action data, including the items clicked and watched, along with the user's explicit feedback (e.g., ratings, feelings). The first step in reflection is for the agent to determine what to reflect on, then we 364 prompt the language model to extract insights and cite the particular records that served as evidence 365 for the insights. The post-interaction reflections are fed back into the episodic memory. 366

367 368

369

4 EXPERIMENTS

3704.1 EXPERIMENTAL SETTINGS371

Datasets. We investigate three real-world datasets: MovieLens-1M Harper & Konstan (2015), Steam Kang & McAuley (2018), and AmazonBook McAuley et al. (2015). They are employed for the initialization of each agent — persona and memory modules, as well as self-consistent persona matching. In order to address privacy concerns, the name and gender are omitted. Moreover, for the sake of generality, we do not utilize user-specific information available in these datasets, relying instead on the personas identified in Phase 1 of SimUSER.

Settings. All agents are powered by the GPT-40-mini version of ChatGPT, except when specified

3	ſ	g
3	7	9

391

392

393

394

395

396

397

398

380		MovieLens			AmazonBook			Steam					
381	Method(1:m)	Accuracy	Precision	Recall	F1 Score	Accuracy	Precision	Recall	F1 Score	Accuracy	Precision	Recall	F1 Score
	RecAgent (1:1)	0.5807	0.6391	0.6035	0.6205	0.6035	0.6539	0.6636	0.6587	0.6267	0.6514	0.6490	0.6499
382	RecAgent (1:2)	0.5498	0.7475	0.5269	0.6178	0.6372	0.6520	0.5511	0.5970	0.6240	0.6788	0.5868	0.6290
000	RecAgent (1:3)	0.5077	0.7396	0.3987	0.5181	0.6144	0.6676	0.4001	0.5003	0.5873	0.6674	0.3488	0.4576
383	RecAgent (1:9)	0.4800	0.7491	0.2168	0.3362	0.6222	0.6641	0.1652	0.2647	0.5995	0.6732	0.1744	0.2772
384	Agent4Rec (1:1)	0.6912	0.7460	0.6914	0.6982	0.7190	0.7276	0.7335	0.7002	0.6892	0.7059	0.7031	0.6786
	Agent4Rec (1:2)	0.6466	0.7602	0.5058	0.5874	0.6842	0.6888	0.5763	0.5850	0.6755	0.7316	0.5371	0.5950
385	Agent4Rec (1:3)	0.6675	0.7623	0.4562	0.5433	0.6707	0.6909	0.4423	0.5098	0.6505	0.7381	0.4446	0.5194
386	Agent4Rec (1:9)	0.6175	0.7753	0.2139	0.3232	0.6617	0.6939	0.2369	0.3183	0.6021	0.7213	0.1901	0.2822
500	SimUSER (1:1)	0.7834	0.7942	0.7455	0.7691	0.8012	0.7865	0.7538	0.7698	0.7809	0.7812	0.7651	0.7731
387	SimUSER (1:2)	0.7533	0.7811	0.5666	0.6566	0.7390	0.7887	0.6659	0.7216	0.7512	0.7956	0.6063	0.6883
	SimUSER (1:3)	0.7412	0.7989	0.5124	0.6263	0.6512	0.7432	0.5542	0.6341	0.7299	0.7959	0.5248	0.6351
388	SimUSER (1:9)	0.6776	0.8096	0.3420	0.4794	0.6388	0.7460	0.3113	0.4384	0.7005	0.7720	0.2651	0.3934

differently, with the number of agents set to 1,000. The prompts and implementation details can be found in the Appendix B. Our analysis consists of two perspectives: micro-level and macro-level. **Baselines** We compare the performance of SimUSER against RecAgent Wang et al. (2023c) and Agent4Rec Zhang et al. (2023), which represent the closest comparable methods for systematic recommender system assessment. When possible, we report the results of RecMind Wang et al. (2023e), an agent-based RS. Some experiments involve two versions of SimUSER: SimUSER(zero) and SimUSER(sim), where SimUSER(sim) agents first interact with the recommender platform grounding interactions and filling their memories, before answering the tasks.

399

400

4.2 BELIEVABLY OF SYNTHETIC USERS

401 A central question is: How accurately do synthetic users replicate the behaviors of real humans? 402 In order to appropriately respond to recommendations, synthetic users must possess a clear under-403 standing of their own preferences. Thereby, we query the agents to classify items based on whether their human counterparts have interacted with them or not. Specifically, each of the 1,000 agents 404 is randomly assigned 20 items. The proportion of items that the user has interacted with but was 405 not used for persona and memory modules initialization to those that the user has not interacted 406 with (represented as $y_{ui} = 0$) is set as 1: m with $m \in \{1, 2, 3, 9\}$. We treat the problem as a 407 binary classification task, taking values between 0 and 1. The results are presented in Table 1, from 408 which we can see that SimUSER agents accurately identify items that are aligned with their own 409 tastes. Across varying levels of distractors, SimUSER exhibits significantly superior performance 410 compared to RecAgent and Agent4Rec, as indicated by paired t-tests at 95% confidence level (p < 411 0.002). Our method achieves an accuracy ranging from 78% to 80%, and recall between 26% and 412 76%, whereas Agent4Rec's accuracy spans from 60% to 71%, and recall varies from 19% to 73%. 413 We attribute SimUSER's improvements to the use of a persona module augmented with granular 414 information (e.g., age, occupation) obtained via persona matching. Notably, Agent4Rec displays 415 a tendency to favor newly encountered items positively, resulting in a higher rate of item approval compared to their human counterparts. SimUSER overcomes this drawback by leveraging the KG 416 memory to retrieve similar items, reducing bias towards unfamiliar items. 417

419 4.3 RATING ITEMS

420

418

A key task when interacting with a recommender system is rating items — predicting the rating 421 that a user would give to a particular item. We compare the results of several LLM-based base-422 lines, RecMind-SI (few-shot) Wang et al. (2023e), RecAgent, and Agent4Rec, along with tradi-423 tional recommendation baselines: Matrix Factorization (MF) Koren et al. (2009) and Attentional 424 Factorization Machines (AFM) Xiao et al. (2017). The results of rating prediction are summarized 425 in Table 2. Across all tasks, we observe that the present study considerably outperforms the other 426 LLM-powered agents: RecAgent, Agent4Rec, and RecMind-SI. Such improvement mainly stems 427 from the advantage that SimUSER has access to the knowledge-graph memory that encapsulates 428 priors about items and their relationships with user interactions. In contrast, Agent4Rec usually has a much higher RMSE, which can be attributed to its hallucination when rating items not embedded 429 within its LLM's weights, particularly niche items. As expected leveraging persona matching, as 430 done in SimUSER(persona), reduces the prediction errors compared to SimUSER(no-persona). An-431 other interesting trend is that incorporating a few steps of simulation always decreases the MAE of Table 2: Performance comparison in rating prediction on MovieLens, AmazonBook, and Steam datasets. The best results of each model are marked in **bold** and the second-best results are marked with underline. The improvement of SimUSER over all baselines on overall performance is statistically significant (measured by student's t-test at p < 0.05).

Methods	Movi	eLens	Amazo	nBook	Ste	am
	RMSE	MAE	RMSE	MAE	RMSE	MAE
MF	1.2142	0.9971	1.2928	0.9879	1.3148	1.0066
AFM	1.1762	0.8723	1.3006	1.1018	1.2763	0.9724
RecAgent	1.1021	0.7632	1.2587	1.1191	1.0766	0.9598
RecMind-SI (few-shot)	1.0651	0.6731	1.2139	0.9434	0.9291	0.6981
Agent4Rec	0.7612	0.7143	0.8788	0.6712	0.7577	0.6880
SimUSER(sim · persona)	0.5341	0.4671	0.5919	0.4562	0.6153	0.5686
SimUSER(zero · w/o persona)	0.6872	0.5981	0.7103	0.6667	0.7089	0.6928
SimUSER(zero · persona)	0.6112	0.5353	0.6698	0.5597	0.6844	0.6392
SimUSER(sim · w/o persona)	0.6012	0.5421	0.6811	0.5605	0.6912	0.6420

the model (SimUSER(sim)). This is because the grounded interactions augment the context during decision-making. For instance, when an agent was asked to rate the movie The City of Lost Children, it retrieved movies such as Johnny Mnemonic and Judge Dredd, all featuring similar adventure and sci-fi themes that it had previously enjoyed. In the absence of simulation, the agent relied solely on its (limited) viewing history from the initial dataset. This demonstrates that agents can refine their own preferences for unrated items through interactions with the simulator.

4.4 RATING DISTRIBUTION



Table 3:	Evaluation	of recon	nmendation	strategies	on
a recomr	nendation ta	ask from	the MovieL	ens datase	t.

	$\overline{P}_{\text{view}}$	$\overline{N}_{\mathrm{like}}$	$\overline{P}_{\text{like}}$	$\overline{N}_{\mathrm{exit}}$	$\overline{S}_{\mathrm{sat}}$
Random	0.298	3.17	0.248	2.87	2.70
Рор	0.401	4.12	0.379	2.89	3.35
MF	0.469	5.98	0.446	3.09	3.67
MultVAE	0.521	5.44	0.458	3.21	3.82
LightGCN	0.552	5.50	0.451	3.34	3.97

Figure 2: Comparison of rating distributions between ground-truth and human proxies.

While aligning individual agent ratings with their human counterparts is crucial, it is also necessary that human proxies replicate real-world user behavior at a macro level. This implies ensuring that the distribution of ratings generated by the agents aligns closely with the distributions observed in the original dataset. Figure 2 presents the rating distribution from the MovieLens-1M dataset and the ratings generated by the agents. These results reveal a high degree of alignment between the simulated and actual rating distributions, with a predominant number of ratings at 4 and a small number of low ratings (1-2). While Agent4Rec assigns fewer 1-2 ratings compared to real users, our approach, by retrieving past interactions from the episodic memory, allows agents to contextu-alize their ratings based on a broader and more consistent understanding of their own preferences. Besides, SimUSER features a more more granular representation of users by incorporating detailed personas (e.g., summary of tastes, personality traits), which yields realistic macro-level behaviors.

RECOMMENDER SYSTEM EVALUATION 4.5

Understanding the efficacy of various recommendation algorithms is crucial for enhancing user satisfaction. By simulating human proxies, we can better predict how users will engage with rec-ommender systems, thereby providing valuable interactive metrics. We also conjecture that users experience higher satisfaction with advanced recommendation strategies compared to random rec-

Table 4: Human-likeness score evaluated by GPT-40 in different recommendation domains. Asterisks (*) denote statistically significant improvements over the best baseline (t-test at p < 0.05).

	MovieLens	AmazonBook	Steam
RecAgent	3.01 ± 0.14	3.14 ± 0.13	2.96 ± 0.17
Agent4Rec	3.04 ± 0.12	3.21 ± 0.14	3.09 ± 0.16
SimUSER(w/o persona)	$3.60 \pm 0.17*$	$3.54 \pm 0.19^*$	$3.58 \pm 0.26*$
SimUSER(persona)	4.22±0.17*	3.98±0.18*	3.97±0.24*

494 495

496

ommendations. We tested several recommendation strategies within our simulator, including most 497 popular (Pop), matrix factorization (MF) Koren et al. (2009), LightGCN He et al. (2020), and Mult-498 VAE Liang et al. (2018), using the MovieLens dataset. Upon exiting, agents rated their satisfaction 499 with the recommendation system on a scale from 1 to 10. Ratings above 3 were considered indica-500 tive of a *like*. The metrics collected included average viewing ratio (P_{view}), average number of likes 501 $(\overline{N}_{\text{like}})$, average ratio of likes $(\overline{P}_{\text{like}})$, average exit page number $(\overline{N}_{\text{exit}})$, and average user satisfaction 502 score (S_{sat}). Table 3 summarizes the performance metrics for the various recommendation strate-503 gies. Agents exhibit higher satisfaction with advanced recommendations compared to random and 504 popularity-based methods, consistent with observed real-life trends. LightGCN outperformed both 505 MF and MultVAE across multiple metrics, confirming its effectiveness in providing personalized 506 recommendations. Namely, LightGCN achieved the highest P_{view} and S_{sat} scores, indicating supe-507 rior user engagement and satisfaction. Overall, these findings highlight the potential of generative 508 agents as a cost-effective technique for simulated A/B testing.

509 510

511 4.6 LLM EVALUATOR

512

513 LLM Evaluators Chiang & Lee (2023) have validated LLM as an evaluator, achieving comparable 514 performance with expert evaluators. Thereby, we use GPT-40 to evaluate the interactions and feedback generated by generative agents. That is, we collected the interactions of the agents and asked 515 GPT-40 to differentiate whether these interactions are from humans or AI-generated. A 5-point 516 Likert scale questionnaire was used. In this scale, a higher score indicates a closer resemblance to 517 human-like responses. Although both methods leverage the semantic awareness of LLMs to achieve 518 near-human-level performance, the results in Table 4 demonstrate that our method significantly out-519 performs Agent4Rec. The memory and persona modules are among the main factors contributing 520 to the faithfulness of our method. We also noticed that letting the agent estimate its tiredness, feel-521 ing and emotion greatly enhances the believability and consistency of its responses. On the other 522 hand, in Agent4Rec, tendencies to [EXIT] the recommender system early and provide inconsistent 523 ratings for similar items — ranging from low to high, contribute to suspicions of AI involvement.

524 525

5 CONCLUSION

526 527 528

In this study, we present a simulation framework for leveraging LLMs as believable user proxies. 529 Our two-phase approach includes persona matching and interactive recommender system assess-530 ment, seeking to align user interactions more closely with real-world user behaviors. Key novelties 531 of our method include the integration of detailed user personas, a knowledge-graph memory mod-532 ule, alongside integrating visual cues to emulate visual-driven reasoning. The present study also 533 introduces a post-interaction reflection mechanism to let the agent discover high-level insights from 534 interactions. We evaluate SimUSER across various recommendation domains, including movies, books, and video games. Results demonstrate closer alignment of our agents with their human 536 counterparts at both micro and macro levels. We further explore the influence of thumbnails on user engagement and the significance of reviews in user decision-making. Experimental findings highlight the potential of LLM-driven simulations in bridging the gap between offline metrics and online 538 evaluation, paving the way for advancements in the development of systematic and cost-efficient evaluation of recommender systems.

540 REFERENCES

547

554

558

559

560

567

568

569

570

573

577

578

579

588

589

- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. Using large language models to simulate
 multiple humans and replicate human subject studies. In *International Conference on Machine Learning*, pp. 337–371. PMLR, 2023.
- Richard C Atkinson and Richard M Shiffrin. Human memory: A proposed system and its control processes. In *Psychology of learning and motivation*, volume 2, pp. 89–195. Elsevier, 1968.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17682–17690, 2024.
- Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937*, 2023.
- Konstantina Christakopoulou, Filip Radlinski, and Katja Hofmann. Towards conversational recommender systems. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 815–824, 2016.
 - Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language models. arXiv preprint arXiv:2302.06692, 2023.
- Wenqi Fan, Zihuai Zhao, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Jiliang Tang, and Qing Li. Recommender systems in the era of large language models (llms). *arXiv preprint arXiv:2307.02046*, 2023.
- Michael Färber, Melissa Coutinho, and Shuzhou Yuan. Biases in scholarly recommender systems:
 impact, prevalence, and mitigation. *Scientometrics*, 128(5):2703–2736, 2023.
 - Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. S3: Social-network simulation system with large language model-empowered agents. *arXiv preprint arXiv:2307.14984*, 2023.
- Xiaoxiao Guo, Hui Wu, Yu Cheng, Steven Rennie, Gerald Tesauro, and Rogerio Feris. Dialog-based
 interactive image retrieval. *Advances in neural information processing systems*, 31, 2018.
- Perttu Hämäläinen, Mikke Tavast, and Anton Kunnari. Evaluating large language models in generating synthetic hci research data: a case study. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1–19, 2023.
 - F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. Acm transactions on interactive intelligent systems (tiis), 5(4):1–19, 2015.
- Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. Lightgcn:
 Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*,
 pp. 639–648, 2020.
- Zhankui He, Zhouhang Xie, Rahul Jha, Harald Steck, Dawen Liang, Yesu Feng, Bodhisattwa Prasad
 Majumder, Nathan Kallus, and Julian McAuley. Large language models as zero-shot conversa tional recommenders. In *Proceedings of the 32nd ACM international conference on information and knowledge management*, pp. 720–730, 2023.
 - Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. Large language models are zero-shot rankers for recommender systems. In *European Conference on Information Retrieval*, pp. 364–381. Springer, 2024.
- Xu Huang, Jianxun Lian, Yuxuan Lei, Jing Yao, Defu Lian, and Xing Xie. Recommender
 ai agent: Integrating large language models for interactive recommendations. *arXiv preprint* arXiv:2308.16505, 2023.

598

601

619

630

- 594 Dietmar Jannach and Michael Jugovac. Measuring the business value of recommender systems. ACM Transactions on Management Information Systems (TMIS), 10(4):1–23, 2019. 596
 - Wang-Cheng Kang and Julian McAuley. Self-attentive sequential recommendation. In 2018 IEEE international conference on data mining (ICDM), pp. 197–206. IEEE, 2018.
- Wang-Cheng Kang, Jianmo Ni, Nikhil Mehta, Maheswaran Sathiamoorthy, Lichan Hong, Ed Chi, 600 and Derek Zhiyuan Cheng. Do llms understand user preferences? evaluating llms on user rating prediction. arXiv preprint arXiv:2305.06474, 2023. 602
- Byungwan Koh and Fuquan Cui. An exploration of the relation between the visual attributes of 603 thumbnails and the view-through of videos: The case of branded video content. Decision Support 604 Systems, 160:113820, 2022. 605
- 606 Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender 607 systems. Computer, 42(8):30-37, 2009.
- 608 Wenqiang Lei, Xiangnan He, Yisong Miao, Qingyun Wu, Richang Hong, Min-Yen Kan, and Tat-609 Seng Chua. Estimation-action-reflection: Towards deep interaction between conversational and 610 recommender systems. In Proceedings of the 13th International Conference on Web Search and 611 Data Mining, pp. 304-312, 2020. 612
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Com-613 municative agents for" mind" exploration of large language model society. Advances in Neural 614 Information Processing Systems, 36:51991–52008, 2023a. 615
- 616 Jinming Li, Wentao Zhang, Tian Wang, Guanglei Xiong, Alan Lu, and Gerard Medioni. Gpt4rec: A 617 generative framework for personalized recommendation and user interests interpretation. arXiv 618 preprint arXiv:2304.03879, 2023b.
- Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang 620 Liu. Agent hospital: A simulacrum of hospital with evolvable medical agents. arXiv preprint 621 arXiv:2405.02957, 2024a. 622
- 623 Yang Li, Kangbo Liu, Ranjan Satapathy, Suhang Wang, and Erik Cambria. Recent developments 624 in recommender systems: A survey. *IEEE Computational Intelligence Magazine*, 19(2):78–95, 2024b. 625
- 626 Dawen Liang, Rahul G Krishnan, Matthew D Hoffman, and Tony Jebara. Variational autoencoders 627 for collaborative filtering. In Proceedings of the 2018 world wide web conference, pp. 689-698, 628 2018. 629
 - Yuqing Liu, Yu Wang, Lichao Sun, and Philip S Yu. Rec-gpt4v: Multimodal recommendation with large vision-language models. arXiv preprint arXiv:2402.08670, 2024.
- 632 Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 633 Learning word vectors for sentiment analysis. In Proceedings of the 49th annual meeting of the 634 association for computational linguistics: Human language technologies, pp. 142–150, 2011.
- 635 Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. Image-based rec-636 ommendations on styles and substitutes. In Proceedings of the 38th international ACM SIGIR 637 conference on research and development in information retrieval, pp. 43-52, 2015. 638
- 639 Tien T Nguyen, F Maxwell Harper, Loren Terveen, and Joseph A Konstan. User personality and user 640 satisfaction with recommender systems. Information Systems Frontiers, 20:1173–1189, 2018.
- 641 Paul Owoicho, Ivan Sekulic, Mohammad Aliannejadi, Jeffrey Dalton, and Fabio Crestani. Ex-642 ploiting simulated user feedback for conversational search: Ranking, rewriting, and beyond. In 643 Proceedings of the 46th International ACM SIGIR Conference on Research and Development in 644 Information Retrieval, pp. 632-642, 2023. 645
- Joon Sung Park, Joseph C O'Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and 646 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. arXiv preprint 647 arXiv:2304.03442, 2023.

684

688

689

- Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 6, 2023.
- Yizhou Sun, Jiawei Han, Xifeng Yan, Philip S Yu, and Tianyi Wu. Pathsim: Meta path-based top-k
 similarity search in heterogeneous information networks. *Proceedings of the VLDB Endowment*, 4(11):992–1003, 2011.
- Yueming Sun and Yi Zhang. Conversational recommender system. In *The 41st international acm sigir conference on research & development in information retrieval*, pp. 235–244, 2018.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. arXiv preprint arXiv:2305.16291, 2023a.
- Hanbing Wang, Xiaorui Liu, Wenqi Fan, Xiangyu Zhao, Venkataramana Kini, Devendra Yadav, Fei
 Wang, Zhen Wen, Jiliang Tang, and Hui Liu. Rethinking large language model architectures for
 sequential recommendations. *arXiv preprint arXiv:2402.09543*, 2024a.
- Jiawei Wang, Renhe Jiang, Chuang Yang, Zengqing Wu, Makoto Onizuka, Ryosuke Shibasaki, and Chuan Xiao. Large language models as urban residents: An llm agent framework for personal mobility generation. *arXiv preprint arXiv:2402.14744*, 2024b.
- Lei Wang, Jingsen Zhang, Xu Chen, Yankai Lin, Ruihua Song, Wayne Xin Zhao, and Ji-Rong
 Wen. Recagent: A novel simulation paradigm for recommender systems. *arXiv preprint arXiv:2306.02552*, 2023b.
- Lei Wang, Jingsen Zhang, Hao Yang, Zhiyuan Chen, Jiakai Tang, Zeyu Zhang, Xu Chen, Yankai
 Lin, Ruihua Song, Wayne Xin Zhao, et al. User behavior simulation with large language model
 based agents. *arXiv preprint arXiv:2306.02552*, 2023c.
- Kiaolei Wang, Xinyu Tang, Wayne Xin Zhao, Jingyuan Wang, and Ji-Rong Wen. Rethinking the evaluation for conversational recommendation in the era of large language models. *arXiv preprint arXiv:2305.13112*, 2023d.
- Kinyuan Wang, Liang Wu, Liangjie Hong, Hao Liu, and Yanjie Fu. Llm-enhanced user-item interactions: Leveraging edge information for optimized recommendations. *arXiv preprint* arXiv:2402.09617, 2024c.
- Yancheng Wang, Ziyan Jiang, Zheng Chen, Fan Yang, Yingxue Zhou, Eunah Cho, Xing Fan, Xiaojiang Huang, Yanbin Lu, and Yingzhen Yang. Recmind: Large language model powered agent for recommendation. *arXiv preprint arXiv:2308.14296*, 2023e.
- Zhenduo Wang, Zhichao Xu, Vivek Srikumar, and Qingyao Ai. An in-depth investigation of user
 response simulation for conversational search. In *Proceedings of the ACM on Web Conference* 2024, pp. 1407–1418, 2024d.
 - Yunjia Xi, Weiwen Liu, Jianghao Lin, Xiaoling Cai, Hong Zhu, Jieming Zhu, Bo Chen, Ruiming Tang, Weinan Zhang, Rui Zhang, et al. Towards open-world recommendation with knowledge augmentation from large language models. arXiv preprint arXiv:2306.10933, 2023.
- Jun Xiao, Hao Ye, Xiangnan He, Hanwang Zhang, Fei Wu, and Tat-Seng Chua. Attentional factorization machines: Learning the weight of feature interactions via attention networks. *arXiv* preprint arXiv:1708.04617, 2017.
- Fan Yang, Zheng Chen, Ziyan Jiang, Eunah Cho, Xiaojiang Huang, and Yanbin Lu. Palr: Personalization aware llms for recommendation. *arXiv preprint arXiv:2305.07622*, 2023a.
- Heng Yang, Chen Zhang, and Ke Li. Pyabsa: A modularized framework for reproducible aspectbased sentiment analysis. pp. 5117–5122, 2023b.
- Se-eun Yoon, Zhankui He, Jessica Maria Echterhoff, and Julian McAuley. Evaluating large language models as generative user simulators for conversational recommendation. *arXiv preprint arXiv:2403.09738*, 2024.

702 703 704	An Zhang, Leheng Sheng, Yuxin Chen, Hao Li, Yang Deng, Xiang Wang, and Tat-Seng Chua. On generative agents in recommendation. <i>arXiv preprint arXiv:2310.10108</i> , 2023.
704 705 706	Jizhi Zhang, Keqin Bao, Wenjie Wang, Yang Zhang, Wentao Shi, Wanhong Xu, Fuli Feng, and Tat-Seng Chua. Prospect personalized recommendation on large language model-based agent
707	platform. arxiv preprint arXiv:2402.18240, 2024a.
708	Junjie Zhang, Yupeng Hou, Ruobing Xie, Wenqi Sun, Julian McAuley, Wayne Xin Zhao, Leyu
710	Lin, and Ji-Kong Wen. Agentci: Collaborative learning with autonomous language agents for recommender systems. In Proceedings of the ACM on Web Conference 2024, pp. 3670–3680
711	2024b.
712	
713	Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. Deep learning based recommender system: A survey and new perspectives. ACM computing surveys (CSUR), 52(1):1–38, 2019.
714	
715	Shuo Zhang and Krisztian Balog. Evaluating conversational recommender systems via user simula-
716	data mining, pp. 1512–1520, 2020.
718	Shuo Zhang, Mu-Chun Wang, and Krisztian Balog. Analyzing and simulating user utterance refor-
719	mulation in conversational recommender systems. In Proceedings of the 45th International ACM
720	SIGIR Conference on Research and Development in Information Retrieval, pp. 133–143, 2022.
721	Kesen Zhao, Shuchang Liu, Qingneng Cai, Xiangyu Zhao, Ziru Liu, Dong Zheng, Peng Jiang, and
723	Kun Gai. Kuaisim: A comprehensive simulator for recommender systems. Advances in Neural
724	Information Processing Systems, 36:44880–44897, 2023.
725	True 7 has Inlien McAulan and Innin King, Languaging as dial source time to income a successful and
726	ranking for collaborative filtering. In <i>Proceedings of the 23rd ACM international conference on</i>
727	conference on information and knowledge management pp 261–270 2014
728	congerence on agermanen and meagermanagement, ppr 201 210, 201 n
729	
730	
731	
732	
733	
734	
735	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
740	
750	
751	
752	
753	
754	
755	

756 A APPENDIX

757

758 **Experimental Settings.** We separate the dataset into training, validation, and test sets (80/10/10%). 759 Relationships between users and items from the training/validation and test sets were excluded 760 from the knowledge graph memory to prevent data leakage. In this paper, we report results 761 for SimUSER with simulation SimUSER(sim), and without simulation SimUSER(zero). In 762 SimUSER(zero), the agent's memory module is initialized from the history of its human counterpart. When the review score for an item is greater than 4, the agent stores a memory entry in 763 the form I liked {item_name} based on my review score of {score}. For a 764 score of 2 or below, the following format is utilized I disliked {item_name} based on 765 my review score of {score}. Neutral scores result in the entry I felt neutral 766 about {item_name} based on my review score of {score}. In SimUSER(sim), 767 agents can also interact with the recommender systems (training set) for up to 20 pages or exit 768 the system at any time. The corresponding interactions are used to enhance the memory module. 769 In all the experiments, items rated ≥ 4 are considered as liked by the user, while items ≤ 2 are 770 considered as disliked. These interactions are stored both as plain text in the episodic memory 771 and as relationships in the knowledge graph memory. These simulated interactions with the RS 772 are stored in the episodic memory with the following format: The recommender system 773 recommended the following {item_type} to me on page {page_number}: {name_all_items}, among them, I selected {watched_items} and rate 774 them {ratings} respectively. I dislike the rest {item_type} items: 775 {dislike_items}. 776

In some sets of experiments, we report performance without persona matching SimUSER(w/o persona), and with persona matching SimUSER(persona). In the absence of persona matching, personality traits, age, occupation and taste summary are omitted from the prompts. Matrix factorization (MF) is utilized as the recommender model unless specified otherwise. In our simulator, agents are presented with four items n = 4 per page and allowed to interact by viewing and rating items based on their preferences. When the output of the LLM deviated from the desired format, resulting in errors, the LLM was re-prompted with the following instruction: You have one more chance to provide the correct answer.

785 As mentioned above, we leverage GPT-40-mini as the LLM backbone in all the experiments unless stated differently. We use $\alpha = 0.8$ to balance item-item similarity with user-item similarity. We 786 set $k_2 = 3$ when retrieving similar items from the knowledge graph-memory, and $k_1 = 5$ for the 787 episodic memory. The titles and ratings of retrieved items from the knowledge graph are concate-788 nated to condition decision-making prompts. Empirically, we set the weight of node embeddings 789 to 0.25 when computing top- k_2 scores. Documents and embedding of text (E) were obtained us-790 ing text-embedding-3-small. Given the average rating \bar{R} of a user: $\bar{R} = \frac{1}{N} \sum_{i=1}^{N} r_{ui}$, the 791 pickiness level $P(\bar{R})$ of a user was determined based on the following thresholds: 792

$$P(\bar{R}) = \begin{cases} P_1 & \text{if } \bar{R} \ge 4.5\\ P_2 & \text{if } 3.5 \le \bar{R} < 4.5\\ P_3 & \text{if } \bar{R} < 3.5 \end{cases}$$

where P_1 corresponds to *not picky*, P_2 corresponds to *moderately picky*, and P_3 corresponds to *extremely picky*.

In Appendix F.6, we compare the results of SimUSER taking as input: 1) the original movie poster,
a random screenshot from the movie trailer on YouTube, 3) the original movie poster distorted with a blue color filter (hue=30, lightness=30, saturation=30).

An illustration of the method is provided in Figure 3, detailing the interaction between its components and their roles within the proposed framework.

805

793 794

796

B ADDITIONAL IMPLEMENTATION DETAILS

806 807

Persona Matching 10 candidate personas are generated and used as prior for self-consistency check.
 To facilitate cost-efficient persona identification, up to 50 items are randomly sampled from the user's viewing history and passed to query the LLM. For simplicity, consistency check also utilizes



Figure 3: The SimUSER framework for evaluating a movie recommender system.

 $\rho = 50$ items and j = 3 for other users \bar{u} . The j = 3 subsets of the user's interactions are compared with those of 30 other users, which are randomly selected. The following occupations were introduced in the persona-matching prompt in order to enhance diversity:

828		
829	Category	Values
830	Occupation	
831	occupation	Other
832		Academic/educator
833		Artist
834		Clerical/admin
835		College/grad student
836		Customer service
837		Doctor/health care
838		Executive/managerial
839		Farmer
840		Homemaker
0.11		K-12 student
041		Lawyer
842		Programmer
843		Retired
844		Sales/marketing
845		Scientist
846		Self-employed
847		Technician/engineer
848		Tradesman/craftsman
849		Unemployed
850		Writer
000		
851		

We defined the following age categories:

Category	Values
Age	
1	Under 18
18	18-24
25	25-34
35	35-44
45	45-49
50	50-55
56	56+

Trait	Range
Openness	(1-3)
Conscientiousness	(1-3)
Extraversion	(1-3)
Agreeableness	(1-3)
Neuroticism	(1-3)

For personality features, we listed the five Big Five personality traits and their possible ranges:

The score of each big-5 facet was then mapped to the corresponding text in: low, medium, high.

C PROMPTS

In this section we describe the prompts used to queried the LLM, with those marked by blue headers constituting the core framework of SimUSER. Prompts highlighted in purple are task-specific prompts utilized in some of the experiments discussed in this paper. Note that for books (Amazon-Book) and video games (Steam), watching / watch were respectively replaced with reading / read and playing / play. For the sake of generality and readability, we report only the watching/watch version.

C.1 PERSONA MATCHING

The following prompt was used to generate personas from historical data:

Persona Generation Prompt

You are an AI assistant specializing in analyzing human preferences and understanding personas. You will analyze the items liked and disliked by an individual.

The task is as follows: I will provide you with a list of items that the individual liked and disliked. Based on this information, you need to predict the most suitable persona for this individual. The persona includes the individual's occupation, age range, and Big5 personality traits (each scored from 1 to 3, where 1 is low, 2 is medium, and 3 is high).

```
** Liked Items: **
```

1. {item 1}

2. {item 2}

...

...

** Disliked Items: ** 1. {item 1} 2. {item 2}

Based on the preferences listed above, predict the most suitable persona for this individual. Ensure that the predicted persona is coherent with the liked and disliked items, accurately reflecting the individual's tastes and preferences.

Response Format:

9121. **Occupation:** [choose from <occupation list>]9132. **Age Range:** [choose from <age list>]9143. **Big5 Personality Traits:**915- **Openness:** [score from 1 to 3]916- **Conscientiousness:** [score from 1 to 3]

- **Extraversion:** [score from 1 to 3]

- **Agreeableness:** [score from 1 to 3] - **Neuroticism:** [score from 1 to 3]

Please provide your prediction based on the provided liked and disliked items, ensuring the persona is coherent with the individual's tastes and preferences.

Scores of the candidate personas were estimated as:

Persona-Matching Prompt

Pretend to be a {candidate_persona}. Here is a list of {item_type} you like and dislike:

** Liked Items: **
1. {item 1}
2. {item 2}
...
** Disliked Items: **
1. {item 1}
2. {item 2}
...
** Disliked Items: **
1. {item 1}
2. {item 2}
...
Here is a summary of your tastes: {taste}
On a scale of 1 to 10, where 1 is the least likely and 10 is the most likely, rate how likely you
are to like and dislike similar {item_type} in the future, based on your own taste. Explain
the reason for your rating.
Response Format:
Rating:
Reason:

where {item_type} refers to the item type, specifically *movie*, *game*, *book* for MovieLens, Steam, and AmazonBook, respectively.

C.2 TASTE SUMMARY PROMPT

The summary of the user's preferences, {taste}, was obtained by querying the LLM with this prompt:

Taste Summary Prompt

Based on the following lists of {item_type} the user liked and disliked, generate a comprehensive summary of the user's taste, preferences, and aversions.

```
** Liked {item_type}:**

1. {item 1}

2. {item 2}

...
```

** Disliked {item_type}:**
1. {item 1}
2. {item 2}

•••

972

973 974

975

976 977

978

979

980 981

982

983

984

985 986

987

988

989

990

991 992

993

994 995 996

997 998

999

1000 1001

1002 1003

1004

1008

1010

1011

1012

1013

1014

1015

1016

1017

1018

Instructions:

1. Analyze the genres, themes, styles, and any other relevant features of the liked {item_type}.

2. Contrast these with the features of the disliked {item_type}.

3. Summarize the user's overall taste, specific preferences, and aversions in a clear and concise manner.

Response Format:

- Summary of User's Taste: [summary]

C.3 CAPTION GENERATION

In the implementation of the perception module, a caption is generated for each item using GPT-4o as the LLM core, with the following request:

Caption Generation Prompt

You are an expert in generating detailed and descriptive captions for images. You are provided with an image representing a {item_type} poster or a product. Your task is to generate a detailed, descriptive, and concise caption for the image. The caption should accurately describe the visual elements and provide context to help someone decide if they would be interested in clicking on it, as if they were seeing it themselves on a website.

Data Formats:

1. **Image**: A single image of a {item_type} poster or a product.

Task:

- Generate a detailed and concise caption that includes: - Visual details: Colors, shapes, objects, actions, and notable features present in the image.

- Contextual information: Scene or item context, genre, mood, characters for {item_type}, or brand, usage, and key features for products.
- Emotional tone: The atmosphere or emotional appeal conveyed by the image.
- Unique selling points: Highlight any unique features that make the item special.
- Clarity and precision: Ensure the description is clear, precise, and accurately represents the image.
- Conciseness: Be detailed yet concise, avoiding overly long descriptions.

Example:

1020The movie poster for "Avatar: The Way of Water" features four prominently displayed blue-skinned1021Na'vi characters with distinct facial expressions, conveying a sense of intensity and determination.
The background showcases a breathtaking sunset over a vast ocean, with rocky islands emerging from
the water. Below the Na'vi, a dynamic scene captures a Na'vi warrior riding a large, winged creature
above the ocean waves, hinting at adventure and freedom. The title "Avatar: The Way of Water" is
boldly displayed at the bottom, along with the release date "December 15" The overall mood of the
poster is epic and visually stunning, promising an immersive and thrilling cinematic experience set in

the richly detailed world of Pandora.

Please generate the caption based on the provided image.

where the image is a thumbnail of the item. Thumbnail images were systematically scrapped, using the product name as the request.

C.4 RETRIEVAL FROM THE EPISODIC MEMORY

When we retrieve documents from the episodic memory, follow-up questions are generated based on the initial {query}, as follows:

Follow-up question Prompt

Instruction:

Given the following query, generate 1-3 follow-up questions that are directly related to the original query. These follow-up questions should help retrieve documents that provide further relevant yet related information on the topic.

Query:

{query}

Follow-up questions:

1049 Note that caching can be employed to reduce the computational footprint of generating follow-up questions

C.5 WATCHING ITEMS

SimUSER agents rely on the following prompt to decide which items to watch. When interacting with the simulator, the watching and rating/feeling prompts are combined in a chain-of-thoughts fashion as a single query to speedup inference:

Watching Item Prompt

System Prompt

You excel at role-playing. Picture yourself as a user exploring a {item_type} recommendation system. You have the following habits:

- Your engagement trait is described as: {engagement}
- Your conformity trait is described as: {conformity}
- Your variety trait is described as: {variety}

The engagement characteristic pertains to the frequency of your {item_type}-watching habits. The conformity characteristic measures the degree to which your ratings are influenced by historical ratings. The variety characteristic gauges your likelihood of watching {item_type} that may not align with your usual taste.

Pretend to be a {persona_description}. You are {pickiness} about {item_type}. You only watch {item_type} that align with your taste.

Recommended List
{item_list}

Instructions

- 1. Review each {item_type} in the ## Recommended List ##.
- 1078
 2. For each {item_type}, decide if it aligns with your taste. If it aligns, decide to watch it and provide a brief reason. If it does not align, provide a brief reason.

3. Rate each {item_type} you choose to watch on a scale of 1 to 5 and provide a brief feeling or reaction. If you do not watch the {item_type}, mark the rating as N/A.
4. Follow the exact format for each {item_type} as shown below.

Format

{item_type}: [item name]; WATCH: [yes or no]; REASON: [brief reason]; RATING: [1 to 5 if watched, N/A if not watched]; FEELING: [brief feeling]

Example

MOVIE: The Lion King (1994); WATCH: yes; REASON: It's an animated movie with a touch of adventure and fantasy, which aligns with my preferences; RATING: 5; FEELING: It was a magical and emotional journey that I could not resist.

MOVIE: Titanic (1997); WATCH: no; REASON: Does not fit my preference for animated movies or fantasy elements; RATING: N/A; FEELING: N/A

Your Response

where {item_type} specifies the type of item, such as *movie*, *game*. The placeholders {engagement}, {conformity}, and {variety} denote habits estimated by the persona module. {persona_description} is a plain text description of the user's persona (see below). {pickiness} reflects the user's pickiness regarding the items. Finally, {item_list} comprises the list of items recommended by the system, which comprises similar items retrieved from the knowledge graph memory.

C.6 RATING ITEM/FEELING

The rating/feeling prompt is derived from the *watching* prompt. The watching prompt was substituted with the following *rating* prompt in experiments that only require to rate the items — when the agents do not have to decide which items to watch:

Rating Item Prompt

System Prompt

You excel at role-playing. Picture yourself as a user exploring a {item_type} recommendation system. You have the following habits:

- Your engagement trait is described as: {engagement}
- Your conformity trait is described as: {conformity}
- Your variety trait is described as: {variety}

The engagement characteristic pertains to the frequency of your {item_type}-watching habits. The conformity characteristic measures the degree to which your ratings are influenced by historical ratings. The variety characteristic gauges your likelihood of watching {item_type} that may not align with your usual taste.

Pretend to be a {persona_description}. You are {pickiness} about {item_type}. You only watch {item_type} that align with your taste.

Recommended List
{item_list}

Instructions

- 1. Review each {item_type} in the ## Recommended List ##.
- 2. For each {item_type}, decide if it aligns with your taste.
- 1132 3. Rate each {item_type} you choose to watch on a scale of 1 to 5 and provide a brief feeling / reaction.

$\left[\right]$	4. Follow the exact format for each {item_type} as shown below.
	<pre>### Format {item_type}: [item name]; RATING: [1 to 5 if watched, N/A if not watched]; FEELING: [brief feeling] ### Example MOVIE: The Lion King (1994); RATING: 5; FEELING: It's an animated movie with a touch of adventure and fantasy, which aligns with my preferences. It was a magical and</pre>
	emotional journey that I could not resist. MOVIE: Titanic (1997); RATING: 2; FEELING: Not particularly interested in romantic- dramatic themed movies. ### Your response
As frc foi	done in the Watching Prompt (Appendix C.5), this prompt aggregates similar movies retrieved m the KG memory to reduce hallucination and enhance the reasoning context (see below <i>page</i> <i>mat</i>). Namely, each presented movie is accompanied with a list of similar movies, identified sed on relationships between users and items, as well as semantic information.
С.	7 Page Format
Гŀ	e pages of the recommender simulator ({item_list}) were formatted as follows:
	Page Format
	PAGE {page_number} <- {item_title} -> <- History ratings: {item_rating} -> <- Summary: {item_description} -> <- Thumbnail: {thumbnail} -> <- Similar items: {similar_items} ->
	<- {item_title} -> <- History ratings: {item_rating} -> <- Summary: {item_description} -> <- Thumbnail: {thumbnail} -> <- Similar items: {similar_items} ->
vł io it	tere {page_number} represents the page number in the recommender simulation, {item_title} de- tes the title or name of the item, {item_rating} is the average historical rating of the item, and em_description} provides a description of the item. By default, a short description is displayed, t if the agent decides to click on the item, {item_description} is substituted with a more detailed scription.
W (2) the Tr ge tv	hen available, we leverage short descriptions provided in the original Agent4Rec paper Hou et al. (24). In other domains, we employ GPT-40 to generate the content description of each item using following prompt: Summarize the {item_type} {title} with one sentence. e answer cannot include the {item_type} title. Detailed descriptions were merated using the following prompt: Summarize the {item_type} {title} with o to three sentences. The answer cannot include the {item_type} tle. The response is used as the item description.
Th ite a t ** the	e {thumbnail} placeholder contains the caption generated by the captioner model based on the m's thumbnail. Note that {thumbnail} is an optional field that can be omitted if the item lacks humbnail. {similar_items} lists similar items retrieved from the knowledge graph memory, with title (rating/5) ** separated by commas as the format. The rating is replaced with e user's individual rating when available, or alternatively, with the average rating of other users

stored in the KG memory.

1187

1186 C.8 PERSONA DESCRIPTION

The {persona_description} that describes in plain text the persona of a user was set as:

1190 1191 {neuroticism} neuroticism. Summary of your tastes: {taste}. 1192 1193 1194 where the big-5 personality traits were replaced with the corresponding level (low, medium, high), 1195 based on the score obtained during the persona matching stage. {taste} represents the user's preferences, summarized by s_u that is obtained during persona matching, see Appendix C.2. 1196 1197 1198 C.9 POST-INTERACTION REFLECTION 1199 The post-interaction prompt was formatted as: 1200 1201 Post-Interaction Prompt 1202 1203 You are an {item_type} recommendation system. Analyze the user's recent interactions 1204 with {item_type} and their feedback to gain insights into their preferences and suggest 1205 future content. 1207 **###** User Recent Interactions: 1208 - Item watched: {item_title} - History rating: {item_rating} 1209 - User feedback: {rating/feeling} 1210 1211 **###** Example Salient Questions: 1212 - What genres does the user consistently rate highly? 1213 - Are there specific themes or elements that the user enjoys? 1214 - How does the user's preference for certain genres compare to their past behavior? 1215 1216 ### Tasks: 1217 1. Identify and list three most salient questions about the user's preferences and behaviors 1218 based on the recent interactions. 2. Provide three high-level insights about the user's current preferences. Cite the particular 1219 records that served as evidence for each insight. 1220 ### Your Response: 1222 1223 where {rating/feeling} consists of the ratings and feelings expressed by the synthetic user. 1224 1225 1226 C.10 NEXT ACTION 1227 The next action is selected by querying the large language model with this prompt: 1228 1229 Action Prompt 1230

Instructions

1231

1232

1233

1234

1236 1237

1239

1240 1241 You have the following habits:

- Your engagement trait is described as: {engagement}
- Your conformity trait is described as: {conformity}
- Your variety trait is described as: {variety}

The engagement characteristic pertains to the frequency of your {item_type}-watching habits. The conformity characteristic measures the degree to which your ratings are influenced by historical ratings. The variety characteristic gauges your likelihood of watching {item_type} that may not align with your usual taste.

Persona Format

1188

1189

{age} user. You are a {occupation}. You have a {openness} openness, {conscientiousness} conscientiousness, {extraversion} extraversion, {agreeableness} agreeableness, and

1. Generate an overall feeling based on your memory, in accordance with your engagement trait and your satisfaction with the recommender system.
- Evaluate the relevance and interest of each previously shown {item_type}. Consider how well the previous recommendations aligned with your tastes and interests.
- If your overall feeling is positive, write: POSITIVE: [reason].
- If it's negative, write: NEGATIVE: [reason].
2. Assess your level of fatigue. Consider that you may become tired more easily if yo have an inactive engagement trait. Indicate your level of tiredness: NOT TIRED, A LITTL TIRED, VERY TIRED.
 Assess your current emotion. Reflect on your past interactions to estimate your curre emotion. Choose one emotion from the following list that best describes how you fee CURIOUS, FRUSTRATED, EXCITED, NEUTRAL, OVERWHELMED. Select the next decision:
- Decide whether to exit browsing based on your overall feeling, engagement trai and tiredness. You will exit the recommender system if you have negative feeling or if you are tired, especially if you have a low engagement trait. To leave, writ [EXIT]; Reason: [brief reason].
- Decide whether to continue browsing based on your overall feeling, engageme trait, and tiredness. To continue browsing, write: [NEXT]; Reason: [brief reason
- Click on an item to get its full description. To click on an item, write: [CLICK Reason: [brief reason].
- Go back to the previous page. To go back, write: [PREVIOUS]; Reason: [brid reason].
 ### Example 1 1. Overall Feeling: POSITIVE: The recommendations align well with my tastes and the system is easy to use. 2. Tiredness Level: A LITTLE TIRED 3. Emotion: EXCITED 4. Decision: [NEXT]; Reason: I am enjoying the recommendations, and feel that I can continue browsing for a bit longer.
 ### Example 2 1. Overall Feeling: NEGATIVE: The recommendations do not match my preferences, and find the system frustrating. 2. Tiredness Level: VERY TIRED 3. Emotion: FRUSTRATED 4. Decision: [EXIT]; Reason: I am dissatisfied with the recommendations and too tired to continue browsing.
 ### Your Turn 1. Overall Feeling: 2. Tiredness Level: 3. Current Emotion: 4. Decision:

• **Continue browsing** ([NEXT]): If agents have positive feelings and sufficient energy, they may decide to continue browsing.

- Click on an item to get its full description ([CLICK]): When agents want more information about a specific item, they may select this option.
- Go back to the previous page ([PREVIOUS]): If agents wish to return to a previous set of recommendations, they may choose this action.

1302 C.11 Post-Interview Prompt

The prompt presented to each agent for post-interview is as follows:

Post-Interview Prompt

How satisfied are you with the recommender system you recently interacted with?

Instructions:

- 1. Rating: Provide a rating from 1 to 10.
- 2. Explanation: Explain the reason for your rating.

Response Format:

- RATING: [integer between 1 and 10]
- REASON: [detailed explanation]

C.12 BELIEVABLY OF SYNTHETIC USER PROMPT

In Section 4.2, the rating prompt is modified with the following instructions:

Believably of Synthetic User Prompt

Instructions

1. Review each {item_type} in the ## Recommended List ##.

2. For each {item_type}, classify if you have already interacted with it ("Interacted") or if you have not ("Not Interacted").

C.13 LLM EVALUATOR PROMPT

The prompt below was employed to distinguish between humans and AI-generated interactions:

LLM Evaluator Prompt

Please evaluate the following interactions of an agent with a recommender system, and determine whether it is generated by a Large Language Model (LLM) AI or a real human: {interaction logs}

Please rate on a scale of 1 to 5, with 1 being most like an AI and 5 being most like a human.

D

- The knowledge graph memory is initially populated with historical data from datasets such as

KNOWLEDGE-GRAPH CONSTRUCTION

1344 The knowledge graph memory is initially populated with historical data from datasets such as
 1345 MovieLens, Steam, and AmazonBook. Missing information were retrieved from IMDb and/or
 1346 Wikipedia.

MovieLens. The MovieLens dataset provides detailed information about users, movies, and ratings.
We use the version (MovieLens-1M) that includes 1 million user ratings. The construction of the knowledge graph from this dataset involves creating nodes for users and movies and defining edges based on their interactions and attributes. Namely, it incorporates various entity types as nodes:

1350 User, Movie, Genre, Director, Actor, Rating, Release Year, Occupation, and Age. The knowledge 1351 graph contains edges that depict relationships between the following relationships: 1352 • Interaction edges: rated-by, rated, viewed-by, viewed, given-to, gives-rating, 1353 1354 • Preference edges: liked-by, liked, disliked-by, disliked, 1355 • Attribute edges: has-occupation, released-in, contains, has-users, belongs-to, directed-by, 1356 directed, features, acts-in, has-age, has-users. 1357 1358 Steam. The Steam dataset comprises 2,567,538 users, 15,474 games, and 7,793,069 English re-1359 views crawled from Steam, a large video game distribution platform. The dataset also includes rich information that might be useful in future work, like users' play hours, pricing information, media 1360 score, category, developer (etc.). The construction of the knowledge graph from this dataset involves 1361 creating nodes for users and games and defining edges based on their interactions and attributes. The 1362 knowledge-graph memory for the Steam network was created as depicted below. The knowledge 1363 graph incorporates various entity types as nodes: User, Game, Genre, Developer, Publisher, Tag, 1364 Review Sentiment, Price, Release Date, Age, and Occupation. The knowledge graph contains edges 1365 that depict relationships between the following relationships: 1367 Interaction edges: reviewed-by, reviewed, viewed-by, viewed, purchased-by, purchased, • Preference edges: liked-by, liked, disliked-by, disliked, 1369 • Attribute edges: developed-by, published-by, belongs-to-genre, tagged-as, has-review-1370 sentiment, has-price, released-on-date, has-age, has-occupation. 1371 1372 AmazonBook. The AmazonBook dataset provides detailed information about users, books, and 1373 reviews. We use the dataset that includes reviews and book details. The construction of the knowl-1374 edge graph from this dataset involves creating nodes for users and books and defining edges based 1375 on their interactions and attributes. The knowledge graph incorporates various entity types as nodes: User Nodes, Book Nodes, Genre, Author, Publisher, Tag, Review Sentiment, Price, Release Date, 1376 Age, and Occupation. The knowledge graph contains edges that depict relationships between the 1377 following relationships: 1378 1379 Interaction edges: reviewed-by, reviewed, viewed-by, viewed, purchased-by, purchased, 1380 • Preference edges: liked-by, liked, disliked-by, disliked, 1381 • Attribute edges: belongs-to-genre, written-by, published-by, tagged-as, has-review-1382 sentiment, has-price, released-on-date, has-age, has-occupation, 1384 User-tag edges: tags, tagged-by, 1385 • Also bought and also viewed edges: also-bought-by, also-bought, also-viewed-by, also-1386 viewed, 1387 • Bought together edges: bought-by, bought-together. 1388 1389 Graph-Aware Retrieval. To balance computational efficiency and the richness of the captured 1390 relationships, we set the maximum length of the relationship-paths to 4. Here are some examples of 1391 meta-paths for the MovieLens dataset that characterize user-user relationships: 1392 • Movie-User-Movie Rate: Movie $\xrightarrow{rated-by}$ User \xrightarrow{rated} Movie 1393 1394 • Movie-User-Movie Viewed: Movie $\xrightarrow{\text{viewed-by}}$ User $\xrightarrow{\text{viewed}}$ Movie 1395 • Movie-User-Movie Liked: Movie $\xrightarrow{\text{liked-by}}$ User $\xrightarrow{\text{liked}}$ Movie • Movie-User-Movie Disliked: Movie $\xrightarrow{\text{disliked-by}}$ User $\xrightarrow{\text{disliked}}$ Movie 1398 • Movie-Genre-Movie: Movie $\xrightarrow{\text{belongs-to}}$ Genre $\xrightarrow{\text{contains}}$ Movie 1399 1400 • Movie-Director-Movie: Movie $\xrightarrow{\text{directed-by}}$ Director $\xrightarrow{\text{directed}}$ Movie 1401 • Movie-Actor-Movie: Movie $\xrightarrow{\text{features}}$ Actor $\xrightarrow{\text{acts-in}}$ Movie 1402 1403 • Movie-User-Rating-Movie: Movie $\xrightarrow{\text{rated-by}}$ User $\xrightarrow{\text{gives-rating}}$ Rating $\xrightarrow{\text{given-to}}$ Movie

• Movie-Release Year-Movie: Movie $\xrightarrow{\text{released-in}}$ Release Year $\xrightarrow{\text{contains}}$ Movie	
1405 rated-by has-occupation ha	s-users
• Movie-Occupation-Movie: Movie User	>
User $\xrightarrow{\text{Mad}}$ Movie	
• Movie-Age-Movie: Movie $\xrightarrow{\text{rated-by}}$ User $\xrightarrow{\text{has-age}}$ Age $\xrightarrow{\text{has-users}}$ User $\xrightarrow{\text{rated}}$ Movie	
 Steam. Here are some examples of meta-paths for the Steam dataset that characterize relation in the KG: 	ıships
• Game-User-Game Review: Game $\xrightarrow{\text{reviewed-by}}$ User $\xrightarrow{\text{reviewed}}$ Game	
 Game-User-Game Viewed: Game viewed-by User viewed Game 1416 	
• Game-User-Game Liked: Game $\xrightarrow{\text{hked-by}}$ User $\xrightarrow{\text{hked}}$ Game	
• Game-User-Game Disliked: Game $\xrightarrow{\text{disliked-by}}$ User $\xrightarrow{\text{disliked}}$ Game	
• Game-Genre-Game: Game $\xrightarrow{\text{belongs-to-genre}}$ Genre $\xrightarrow{\text{contains}}$ Game	
• Game-Developer-Game : Game $\xrightarrow{\text{developed-by}}$ Developer $\xrightarrow{\text{developed}}$ Game	
• Game-Publisher-Game : Game $\xrightarrow{\text{published-by}}$ Publisher $\xrightarrow{\text{published}}$ Game	
• Game-User-ReviewSentiment-Game: Game $\xrightarrow{\text{reviewed-by}}$ User $\xrightarrow{\text{has-review-set}}$.timent →
1426 ReviewSentiment $\xrightarrow{\text{reviewed}}$ Game	
1427 • Como Price Como: Como has-price Drice Contains	
1428 released on-date contains	
• Game-ReleaseDate-Game: Game $\xrightarrow{\text{Indext-on-out}}$ ReleaseDate $\xrightarrow{\text{Contains}}$ Game	
• Game-Tag-Game: Game $\xrightarrow{\text{tagged-as}}$ Tag $\xrightarrow{\text{contains}}$ Game	
1431 reviewed-by has-occupation has	s-users
4432 • Game-Occupation-Game: Game \longrightarrow User \longrightarrow Occupation $-$	
$\begin{array}{ccc} \text{User} & \xrightarrow{\text{Initial}} & \text{Game} \\ \text{IA2A} & & & & \\ \end{array}$	
• Game-Age-Game: Game $\xrightarrow{\text{reviewed-by}}$ User $\xrightarrow{\text{has-age}}$ Age $\xrightarrow{\text{has-users}}$ User $\xrightarrow{\text{reviewed}}$ Game	
 AmazonBook. The following relationship paths may be considered when retrieving from the zonBook network: 	Ama-
• Book-User-Book Review : Book $\xrightarrow{\text{reviewed-by}}$ User $\xrightarrow{\text{reviewed}}$ Book	
1440 viewed-by Viewed Deck	
• BOOK-USEF-BOOK VIEWED: BOOK — USEF — BOOK	
• Book-User-Book Liked : Book $\xrightarrow{\text{Interm}}$ User $\xrightarrow{\text{Interm}}$ Book	
• Book-User-Book Disliked : Book $\xrightarrow{\text{disliked-by}}$ User $\xrightarrow{\text{disliked}}$ Book	
• Book-Genre-Book : Book $\xrightarrow{\text{belongs-to-genre}}$ Genre $\xrightarrow{\text{contains}}$ Book	
• Book-Author-Book : Book $\xrightarrow{\text{written-by}}$ Author $\xrightarrow{\text{wrote}}$ Book	
• Book-Publisher-Book : Book $\xrightarrow{\text{published-by}}$ Publisher $\xrightarrow{\text{published}}$ Book	
1450 Poole User Devices Confirment Deale Deale reviewed-by User has-review-set	ntiment
$\begin{array}{cccc} \bullet & Dook-User-ReviewSentiment-Book: & Book & \longrightarrow & User & \longrightarrow \\ \text{ReviewSentiment} & \xrightarrow{\text{reviewed}} & Book & & & & \\ \end{array}$	
1452 • Book-Price-Book : Book $\xrightarrow{\text{has-price}}$ Price $\xrightarrow{\text{contains}}$ Book	
1454 Real Palaces Date Back, Date Relaxed-on-date Date Contains Date	
• BOOK-KeleaseDate-BOOK: BOOK → KeleaseDate → Book	
1456 • Book-Tag-Book : Book $\xrightarrow{\text{tagged-as}}$ Tag $\xrightarrow{\text{contains}}$ Book	
• Book-User-Book Bought Together: Book $\xrightarrow{\text{bought-by}}$ User $\xrightarrow{\text{bought-together}}$ Book	

	• Book Hear Book Alea Bought, Dook also-bought-by Hear also-bought, Dook
	• DOOK-USET-DOOK AISO DOUGHT : DOOK
	• Book-User-Book Also Viewed: Book → User → Book
	Reals Occupation Books Deals reviewed-by has-occupation As-users
	• BOOK-Occupation-BOOK: BOOK
	User $\xrightarrow{\text{normal}}$ Book
	• Book-Age-Book : Book $\xrightarrow{\text{reviewed-by}}$ User $\xrightarrow{\text{has-age}}$ Age $\xrightarrow{\text{has-users}}$ User $\xrightarrow{\text{reviewed}}$ Book
Е	Pseudo-Code
We	present the pseudo-code for SimUSER agent.
Ale	corithm 1 SimUSER Algorithm
1.	Input: Historical data <i>H</i> for user <i>u</i>
1. 2.	Output: Simulated interactions and feedback
2. 3.	Phase 1: Persona Matching
$\Delta \cdot$	$\mathcal{P} \leftarrow \text{Generate persona from } H$
	$n \leftarrow \text{Identify best persona} \in \mathcal{P}$ using self-consistency score
6.	Phase 2: Simulate Interactions
7.	Initialize memory module from H.
۶۰ ۲	reneat
٥. ٩٠	Perceive the page and items \triangleright Generate captions
10·	Retrieve similar items from the KG memory
11:	Decide what items to watch
12:	Rate the items and provide feelings
13:	Decide next action a based on satisfaction, fatigue, and emotion
14:	Perform post-interaction reflection
15:	Update memory module
16:	if $a = [EXIT]$ then
17:	break
18:	else
19:	Perform action a
20:	until Maximum number of pages reached
0.1	Return Simulated interactions metrics and feedback

1496 F.1 ALIGNMENT: RATING VS FEELING

1497 Expressing aligned reviews and ratings is of primary importance to simulate realistic human prox-1498 ies. Thus, in this section we delve into the alignment between ratings and sentiments. In detail, 1499 we prompt the agent to assume one has interacted with a certain item, and ask about its rating 1500 and feelings on it. Reviews and ratings from IMDB Maas et al. (2011) are used as ground truth 1501 since MovieLens does not contain reviews. After getting a collection of responses, we conduct 1502 sentiment-based analysis with PyABSA Yang et al. (2023b). We compare the rating and sentiment 1503 distributions of: humans, RecAgent, Agent4Rec, and SimUSER. As depicted in Figure 4, our agents 1504 generate ratings aligned with their opinions. For instance, ratings ≥ 4 are typically associated with 1505 positive sentiments. In contrast, Agent4Rec exhibits a bias towards positive opinions, resulting in 1506 more positive feelings about the items, including when generating low ratings. It is noteworthy that SimUSER agents and genuine humans express similar sentiments at a macro level. 1507

1508

1510

- 1509 F.2 PREFERENCE COHERENCE
- Under this scenario, we aim to evaluate whether agents prefer positive recommendations based on a query. Namely, for each request in the Reddit dataset He et al. (2023), we sample: (1) a comment



Figure 5: Preference coherence (accept/reject task). 'I' stands for incoherent; 'C' stands for coherent.



Figure 6: Distribution of interaction numbers (top) and average ratings (bottom) for 3 groups of personas. The left column does not use persona module while the right employs a persona module.

1566

from this request (positive recommendation) (2) a random comment (negative recommendation). The agent is then prompted to decide which items to *watch*. Ideally, synthetic users should watch only positive recommendations and decline negative ones. Behavior is incoherent when the simulator accepts a negative recommendation. We clearly see in Figure 5 that our agents are overall coherent, but sometimes prefer negative recommendations, its proportion being around 4%. Particularly, Agent4Rec agents often accept recommendations that are not aligned with their age and personality. Upon looking at the responses, we notice that providing explanations, as stated in our prompts, reduces incoherent behaviors in synthetic humans.

1604

1605 F.3 IMPACT OF PERSONA ON USER BEHAVIORS

1607 In real life, user behaviors are driven by factors beyond mere individual tastes, including personal-1608 ity traits and demographic attributes such as age, and occupation. To account for these factors, we 1609 developed a persona module that incorporates these characteristics. Using the MovieLens dataset, 1610 agents were categorized based on their age, occupation (student, office worker, other), and person-1611 ality traits — categorical features were encoded using one-hot representation, via k-means (k=3). We hypothesize that agents with similar preferences may still exhibit distinct behavior patterns due 1612 to their unique persona. Figure 6 depicts the distribution of interaction numbers and average ratings 1613 among agents with different personas. 1614

As shown in the figure, it is evident that agents from different groups demonstrate distinct behaviors. We particularly noticed a higher average rating in group 2 and group 3 than in group 1. Instead, synthetic users without personas display more homogeneous behaviors. Upon looking at the users in each group, our findings indicate that younger users (predominantly in group 3, avg age=19.6) tend to engage more frequently with the system, showing higher interaction numbers compared to older users (group 1, avg age=40.2). This increased activity is likely due to their higher adaptability and interest in exploring new recommendations. In contrast, *older users* exhibit fewer interactions, possibly due to ingrained preferences and less willingness to engage with new content.

F.4 POST-INTERACTION INTERVIEWS

Case Study of Satisfaction Interview

Question: How satisfied are you with the recommender system? Please rate your satisfaction on a scale from 1 to 10 and provide an explanation for your rating.

Rating: 7
Reason: As a 35-year-old film enthusiast with a preference for independent films and documentaries, I found the recommender system fairly effective in suggesting thought-provoking content. It recommended several documentaries and indie films that I enjoyed.

However, the system often leaned towards mainstream and popular films rather than exploring unique and lesser-known titles. This tendency occasionally led to less satisfying recommendations, as I value discovering independent, high-quality films.

Additionally, while the system recognized my interest in dramas and thrillers, it missed specific subgenres like neo-noir and psychological thrillers that I particularly enjoy. Improving the system's ability to identify these nuanced preferences could enhance its effectiveness.

Overall, the system shows strengths in general recommendations but would benefit from a better balance between popular and niche films.

1639 1640

1623

1624 1625

1626 1627

1628

1629

1630

1633

1634

1635

1637

1638

Figure 7: One case study of post-exit feeling interview. The blue color highlights the user's satisfaction score.

1643

1644 A central advantage of LLM-based evaluation is their capacity to generate human-comprehensible 1645 explanations. Such explanations provide valuable insights that can be employed to refine the rec-1646 ommender system. Figure 7 illustrates a post-interaction interview scenario, where the agent rates its experience with the recommender and returns an explanation based on its personal tastes, habits, 1647 emotional memory, and interaction history. The agent noticed that the recommender system gen-1648 erally aligned well with its tastes, particularly recommending movies suitable for its age and per-1649 sonality. However, some factors reduce the overall satisfaction. For instance, despite the agent's 1650 preference for sub-genres like neo-noir and psychological thrillers such as "The Sixth Sense" and 1651 "Apt Pupil", the system frequently recommended mainstream blockbusters like "Star Wars: Episode 1652 I", "The Phantom Menace". These findings underscore the importance of understanding the nuanced 1653 behavior of different recommendation algorithms and their suitability for each group of users.

1654 1655 1656

F.5 RATING ITEMS UNDER HALLUCINATION

1657 In this section, we specifically target items that are unfamiliar to the LLM, seeking to evaluate the 1658 ability of SimUSER to mitigate hallucination through its memory module. Similarly to Section 4.3, users are asked to rate movies (MovieLens). Nevertheless, we exclusively include items that are 1659 detected as unknown to the LLM. These items i are identified by querying the LLM to classify each 1660 movie into one of 18 genres. If the LLM's genre classification matches the actual category g_i , it indicates that the LLM is familiar with the item, and such movies are excluded from the experiment. 1662 From Figure 8, it is evident that while the RMSE values for all methods increase under halluci-1663 nation, the performance degradation of SimUSER is less severe compared to the baselines. This 1664 relative robustness of SimUSER can be attributed to its KG memory, which effectively mitigates 1665 the impact of hallucination by leveraging relationships between users/movies/ratings from previous 1666 interactions. By comparing the unfamiliar movie with these similar, well-known ones, the LLM can anchor its predictions in familiar contexts, reducing the likelihood of hallucinations and leading to 1668 more accurate ratings. Furthermore, beyond mere text-based descriptions, the item captions provide 1669 an additional signal when deciding which items the user likes.

1670

1672

1671 F.6 THUMBNAIL QUALITY EFFECT

1673 Emotions largely shape decision-making in the recommendation domain. At the center of emotion, images are powerful stimuli that motivate our choices. In light of this, a question arises: Can



4.02 4.17 Action 3.70 3.65 3.53 Adventure Animation 3.82 3.74 3.68 Children 3.54 3.41 3.87 3.83 3.74 Comedy Crime 4.12 4.07 Documentary 3.83 3.75 3.66 4.15 4.08 3.97 Movie Genre Drama 3.62 Fantasy 3.75 3.54 4.01 3.93 Film-Noir 3.82 Horro 3,96 4.6 3.76 3.69 3.61 Musica 3.8 Mystery 3.91 3.86 3.72 3.76 3.66 3.55 Romance Sci-Fi 3.90 4.08 3.99 3.6 3.92 Thrille 3.80 3.72 Wai 4.03 4.12 3.4 Western 3.79 3.68 3.54 20 Pages 50 Pages 5 Pages Pages Scrolled

Figure 10: Heatmap showing the impact of biased recommendations on genre ratings over time — exposure effect. The genres and their ratings are displayed after 5, 20, and 50 pages scrolled.

1747

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739

1741

1742

1743

1748

1749 SimUSER be useful in assessing the quality of items' thumbnails? To understand the factors influ-1750 encing ratings, we randomly selected 100 movies from the MovieLens dataset and ask 100 agents 1751 whether they want to watch it. For each movie, we collected three versions of the thumbnails: 1) the 1752 original movie poster, 2) a random screenshot from the movie trailer on YouTube, and 3) the orig-1753 inal movie poster distorted with a blue color filter. Based on the click rates shown in Figure 9, we notice that high-quality thumbnails — original posters, significantly influence users' inclination to 1754 watch a movie. Specifically, original posters lead to higher engagement compared to random screen-1755 shots and color-distorted posters. This result highlights SimUSER's capability to reflect the quality 1756 of item images in decision-making processes, mirroring trends commonly observed in real-world 1757 recommender systems. 1758

1759

1760 F.7 EXPOSURE EFFECT IN RECOMMENDATION

1761 To assess how biased recommendations shape user preferences over time, we introduce a scenario 1762 where the recommender system only recommends two movie categories: *action* and *horror*. It 1763 emulates an exposure effect Färber et al. (2023), where repeated exposures to a particular stimu-1764 lus increase an individual's preference for that stimulus. In the context of recommender systems, 1765 repeated exposure to specific genres could amplify user preferences for those genres. Under this 1766 scenario, we record the average movie ratings for each category after 5, 20, and 50 pages scrolled 1767 by the agents. Namely, the 50 agents are prompted to rate 500 randomly selected movies. Figure 10 1768 illustrates a tendency of the agents to rate items of categories that are over-represented higher during 1769 the interactions with the recommender system, particularly after more than 20 pages. Conversely, 1770 categories that differ significantly from *action* and *horror* genres generally tend to receive lower average ratings. Experimental results validate SimUSER's capability to replicate the exposure effect, 1771 although further research and validation are required with alternative datasets. 1772

1773

1775

1774 F.8 USER REVIEW INFLUENCE

User proxies may help researchers in identifying the psychological effect of reviews on human interactions. To investigate this effect, we modified the recommendation simulator to display a) the number of reviews, b) one random negative comment, or c) one random positive comment for each item on the recommendation page. We report in Table 5 the average viewing ratio \overline{P}_{view} and ratio of likes \overline{P}_{like} . We can see that displaying the number of reviews slightly improves the viewing ratio, especially for items having enough reviews (i.e., more than 20 reviews). This aligns with humans' inclination to select popular items in real-life scenarios. On the other, there is no significant

1783							
1784		Μ	IF	Mult	VAE	Light	GCN
1785 1786	Condition	$\overline{P}_{\text{view}}$	$\overline{P}_{\text{like}}$	$\overline{P}_{\text{view}}$	$\overline{P}_{\text{like}}$	$\overline{P}_{\text{view}}$	$\overline{P}_{\text{like}}$
1787 1788 1789	Origin + With # Reviews + With Negative + With Positive	0.469 0.481 0.428 0.467	0.446 0.479 0.420 0.487	0.501 0.537 0.482 0.553	0.458 0.470 0.435 0.487	0.552 0.549 0.512 0.568	0.451 0.486 0.419 0.486

Table 5: Impact of user reviews on recommender System performance.

1792 Table 6: nDCG@k (k=10) and F1-score@k (k=10) for three recommender systems, using either 1793 offline or SimUSER-generated interactions.

	nDCG@10		F1-score@10		
Method	Offline	SimUSER	Offline	SimUSER	
MF	0.226	0.204	0.165	0.150	
MultVAE	0.288	0.257	0.180	0.172	
LightGCN	0.423	0.454	0.227	0.246	

difference in $\overline{P}_{\text{like}}$ (t-test p > 0.05). Another observation is that displaying negative reviews has a 1803 stronger impact on user behavior than showing positive reviews, with a decrease in both the average viewing ratio and number of likes. One possible explanation is that negative reviews discourage 1805 users from watching an item, while positive reviews primarily encourage users who are already inclined to watch it to proceed with their choice. 1807

1808

1810

1782

1791

1801

F.9 SIMUSER VS. OFFLINE METRICS 1809

We aim to investigate whether SimUSER can reliably estimate traditional metrics such as nDCG@k 1811 (k=10) and F1-score@k (k=10) by comparing the results from traditional offline evaluation with 1812 those from SimUSER-generated interactions. For this purpose, we evaluate three recommender 1813 systems using the MovieLens dataset under identical conditions for both offline and SimUSER-1814 based evaluations. 1815

1816 Table 6 reports the nDCG@k and F1-score@k values (k=10) for both evaluation strategies. In the SimUSER scenario, interactions are generated by our synthetic users, where liked and disliked items 1817 replace the ground-truth interactions from the offline dataset. The results indicate minimal differ-1818 ences between the SimUSER-generated data and the real-world dataset. These differences can be 1819 attributed to real-world scenarios in which users are not constrained by the number of pages and 1820 the frequency of interactions with the recommender system. Interestingly, the results are consis-1821 tent across different systems, with the ranking of models remaining the same between offline and 1822 SimUSER-generated metrics. Overall, these experimental results highlight that SimUSER can reliably measure traditional metrics while also enabling the exploration of system performance across 1824 different user demographic groups, various website settings (e.g., number of items per page), and 1825 different settings of the recommender system.

1826

1827 F.10 ABLATION STUDIES 1828

1829 1830

F.10.1 IMPACT OF THE KNOWLEDGE-GRAPH MEMORY ON SIMUSER

1831 Here, our focus is on evaluating the impact of incorporating a knowledge-graph memory on the performance. Specifically, the goal is to determine whether employing the KG memory, which simulates external influences such as reviews, enhances believability in human proxies. All models follow the same settings as in Sec 4.3. Table 7, highlights that leveraging the KG structure signifi-1834 cantly reduces both RMSE and MAE across different datasets. This module mirrors how our prior 1835 expectations of an item can shape and bias our assessment of it.

Table 7: Performance comparison in rating prediction for agents equipped with (top two rows ♥) and without a KG memory (bottom two rows ♣). Asterisks (*) denote statistically significant improvements when the KG memory is used.

Methods	MovieLens		AmazonBook		Steam	
	RMSE	MAE	RMSE	MAE	RMSE	MAE
SimUSER(zero) 🔍	0.6112*	0.5353*	0.6698	0.5597*	0.6844*	0.6392*
SimUSER(zero) 🐥	0.6712	0.6400	0.6943	0.6220	0.7203	0.6772
SimUSER(sim) 🔍	0.5341*	0.4671*	0.5919*	0.4562*	0.6153*	0.5686*
SimUSER(sim) 🐥	0.6387	0.6551	0.6295	0.4996	0.6654	0.6481

Table 8: Performance of Persona Matching in Predicting Age and Occupation Using the MovieLens-1M Dataset.

Metric	Age	Occupation
Accuracy	0.7230	0.6764
Precision	0.7586	0.6933
Recall	0.7921	0.7430
F1 Score	0.7749	0.7172

1855 1856 1857

1848

1849 1850 1851

1858 F.10.2 PERSONA MATCHING: AGE, OCCUPATION

1859 In this study, we postulate that personas are crucial for capturing the heterogeneity and diversity 1860 present in real-world recommender networks. These attributes significantly shape individual behav-1861 iors and preferences, which subsequently influence the overall dynamics of the system. To evaluate 1862 the effectiveness of our self-consistent persona-matching technique, we conducted an experiment 1863 using the MovieLens-1M dataset. The goal was to predict the age and occupation of users based 1864 on their historical data. This task was formulated as a classification problem. Our results are summarized in table 8. We observe a high degree of alignment between the predicted and actual user 1865 personas, highlighting the effectiveness of Phase 1 in SimUSER. Overall, persona matching turns 1866 out to be reasonably robust for enriching simulated agents with detailed backgrounds, including 1867 domains where explicit demographic data is not readily provided. 1868

1869

1870 F.10.3 PERSONA MATCHING: PERSONALITY

1871 In order to assess the quality of persona matching in predicting personality traits from historical 1872 interaction data, we conduct an additional experiment using the Personality 2018 dataset Nguyen 1873 et al. (2018). The primary objective is to evaluate whether our model could accurately infer users' 1874 Big Five personality traits — Openness, Conscientiousness, Extraversion, Agreeableness, and Neu-1875 roticism, based solely on users' watching history. For a fair comparison, the personality traits within the dataset, as well as the predictions, are normalized to a scale ranging from 0 to 1. We report the 1876 results for various lengths of movie history $\rho \in \{10, 20, 50\}$. This task is formulated as a regression 1877 problem, with the model's performance evaluated using the Mean Absolute Error (MAE) for each 1878 personality trait prediction. Figure 11 summarizes the results, showing that our model achieved an 1879 average MAE of 0.155 across all traits. Besides, the results reveal that using a history length of 50 1880 items reduces the average MAE from 0.279 (10 items) to 0.155, demonstrating that self-consistent 1881 persona matching can reasonably predict personality traits of users from their past interactions.

1882

1883 1884 F.10.4 CHOICE OF FOUNDATION MODEL

1885 We seek to evaluate the performance of our methodology using various foundation models on the 1886 movie rating task. Specifically, we compare the results obtained by employing GPT-4o-mini, GPT-40, Mistral-7b Instruct, Llama-3 Instruct, and Phi-3-mini as the underlying LLMs. The results, 1888 presented in Table 9, demonstrate that the performance of SimUSER is generally robust across 1889 different foundation models. While GPT-40 exhibits significantly lower mean RMSE and MAE (ttest p < 0.05), GPT-40-mini achieves similar performance but with a lower inference time. Mistral-







Figure 12: Human-likeness score evaluated by GPT-40 for SimUSER trained without and with selfask document retrieval.

1966

1968

1944

1945

1946

1947 1948

1949

1950

1951

1952

1953 1954

1957

1958

1959

7b Instruct also performs reasonably well on the MovieLens dataset. On the other hand, Llama-3Instruct and Phi-3-mini, while competitive, show higher errors.

1967 F.10.5 EPISODIC MEMORY ABLATION

In this study, we focus on the episodic memory. We postulate that the proposed self-ask strategy for retrieving documents contributes to the believability of synthetic users. We evaluate agents that use a plain memory, retrieving documents based on the embedding of the query, against agents employing the self-ask strategy. An LLM evaluator is employed to measure the believability of generated interactions. Across all domains, the agents employing self-ask retrieval are more realistic than their simple counterparts (Figure 12), emphasizing the advantages of our episodic memory architecture. This believability gain is largely attributed to the retrieval of more diverse yet relevant "context" during decision-making.

1976

1977 F.10.6 IMPACT OF PERCEPTION MODULE

Finally, we investigate the impact of the perception module on the believability of SimUSER agents.
From Table 10, we can draw several conclusions. Agents consistently exhibit more realistic behavior
when employing a perception module (•), likely due to the inclusion of visual details and unique
selling points associated with the items. Another observation is that the believability gain is lower
on AmazonBook than other datasets. This may be attributed to the nature of books as products —
users are less likely to judge a book based on its cover and more inclined to make decisions based
on its description.

By looking closely at the interactions, we also noticed that agents with behaviors of users featuring different personas are significantly driven by emotional tones provided in the captions. For instance, an agent with high openness may be more inclined to select movies with captions that use positive language like "exciting" or "inspiring". We also postulate that SimUSER (\blacklozenge) may inherit biases from the LLM's interpretation of item descriptions, although these biases could be partially mitigated through adding factual information based on the captions. This suggests that the perception module not only enhances realism at a macro level but also contributes to more visually and emotionally driven engagement.

1993

G ADDITIONAL BASELINE INFORMATION

1995 1996

1997 In this section, we present a comparative analysis of our method against two widely recognized baselines for simulating user interactions: RecAgent and Agent4Rec. Additionally, we examine

1998 Table 10: Human-likeness score evaluated by GPT-40 for SimUSER without () and with () perception module. Asterisks (*) denote statistically significant improvements when the perception 2000 module is activated.

	MovieLens	AmazonBook	Steam
RecAgent	3.01 ± 0.14	3.14 ± 0.13	2.96 ± 0.17
Agent4Rec	3.04 ± 0.12	3.21 ± 0.14	3.09 ± 0.16
SimUSER (♦)	4.12 ± 0.15	3.95 ± 0.16	3.88 ± 0.19
SimUSER (4.22±0.17*	3.98±0.18*	3.97±0.24*

2009 the features of RecMind-SI, which is capable of estimating item ratings. Table 11 offers a detailed 2010 comparison of these methods with human users.

2012 Table 11: Comparison of our method with prior approaches. The 'Emotion/Mood' column rep-2013 resents the agent's capability to incorporate emotional and mood-based reasoning, while 'External Factors' reflects the model's ability to make decisions informed by its pre-existing assumptions 2014 about items. 2015

Name	Interactive	Memory	Image Perception	Emotion/Mood	External Factors	Affordal
RecMind-S	SI X	Personalized + World Knowledge	X	×	1	1
RecAgent	1	Short Term + Long Term	X	X	X	1
Agent4Rec	· /	Episodic	X	1	×	1
SimUSER	1	Knolwedge Graph + Episodic	1	1	1	1
Human	1	Brain	1	1	1	X

Note that, in our analysis, we consider human users as not affordable when compared to the significantly lower cost of leveraging synthetic agents.

Η DISCUSSION

2006 2007 2008

2011

2023

2024

2025 2026

2027

2028 Our study verifies the possibility of synthetic users in recommendation. SimUSER has demon-2029 strated its effectiveness to generate faithful users across a wide range of recommendation domains, 2030 including movies, books, and games. However, we also acknowledge that our method has certain 2031 limitations. 2032

Observed behaviors are well-aligned with existing theories and common behaviors in the recommen-2033 dation domain. Phenomena at micro-level (rating, watching) are manifestations of agent endogenous 2034 behaviors. But why agents possess these behaviors are unexplored due to the black-box nature of the 2035 large language models we adopted. A potential reason could be that LLMs are trained on a massive 2036 corpus that includes texts from various domains. 2037

This paper introduces a strategy to capture external factors in recommendation via a knowledge 2038 graph memory. The key is for the system to retrieve similar items based on the structure of the 2039 graph and their semantic similarities. But in future work, we aim to explore the full potential for 2040 knowledge-graph enriched prompts, such as by dynamically breaking down a complex problem into 2041 smaller subproblems Besta et al. (2024) and retrieving relevant information from the graph using 2042 LLM-generated queries. 2043

The episodic module relies on a self-ask strategy to retrieve documents. Although this technique 2044 introduces additional overhead, we found that it also enhances the diversity and quality of retrieved 2045 documents by extrapolating documents relevant to the query, which is of critical importance in 2046 recommendation systems, where granular user preferences depend on accessing a comprehensive 2047 range of relevant information. 2048

2049 LLMs may replicate biases prevalent in social spaces, such as some groups of individuals being underrepresented. This is problematic if it causes designers to then underlook these peoples' needs 2050 when designing a recommender system. In our experiments, we mitigated this risk by ensuring a 2051 broad range of personas via diverse potential occupations, age, and personalities. We also measured

200-		
2055	Methods	Average Running Time (h)
2056	A gent 4 Rec()	03
2057	Agent/ $\mathbf{Rec}(\mathbf{A})$	0.51
2058	Agent+Ree(+)	0.51
2059	SimUSER(♥)	9.8
2060	SimUSER(+)	0.58
2061		

Table 12: Average running time of Agent4Rec and SimUSER without (♥) and with(♣) parallelization of LLM calls.

2061 2062

the discrepancy between identified and real personas. Our future investigation will focus on ana lyzing underrepresented user groups, as well as evaluating persona matching on a wider range of
 domains (e.g., food).

2066 In terms of evaluation, the assessment of synthetic users' behavior in this study was aimed at a range 2067 of tasks like rating items, measuring the alignment between their feelings and ratings, and LLM-2068 based believability. While these tasks are relevant and provide insight into the model's performance, 2069 they may not fully capture the complexity of real-world user interactions. Real users engage with systems in multifaceted ways, often influenced by contextual factors, evolving preferences, and 2070 interactions that extend beyond simple task execution. To address this limitation, future work could 2071 involve direct comparisons with real user performance, particularly at a macro level. We also seek 2072 to include human-based evaluations to assess the believability of user interactions. 2073

Finally, UX and UI drive our choices and actions in real-world applications. Our simulation, on the other hand, does not fully replicate all those intricate factors, which introduces a gap between real life and simulation. An important future direction is developing an image-based simulator to better capture the complex nature of user experience.

2078

2079 2080 I COST ANALYSIS

2081

We report here the cost of running SimUSER in terms of "cost per 1000 users". Given that the cost per agent may fluctuate slightly due to variations in the number of interactions and the outputs generated by the LLM, we focus on the total cost per 1000 users. Note that the cost per agent rises approximately linearly with the number of agents. Although we express this cost in USD, it may vary significantly depending on the choice of LLM. For our implementation, we utilized GPT-40mini from OpenAI as the underlying model.

As a result, SimUSER achieves a rather low cost of approximately \$12 (\$0.0012/User). A simulation involving 1,000 Agent4Rec agents would cost approximately \$10 (\$0.0010/User). The cost difference mainly stems from the integration of images to enable visual-driven reasoning. Overall, both methods have similar costs and can be employed as *affordable* proxies for human interaction, particularly when compared to their human counterparts. Even though our method incurs a slightly higher cost compared to Agent4Rec, we believe the enhanced believability and greater depth of interactions it offers are key advantages compared to Agent4Rec.

2094 2095 2096

J RUNNING TIME ANALYSIS

2097

2098 In this section, we provide a running time comparison between the proposed method and Agent4Rec. 2099 The simulation involved 1,000 users interacting with the movie recommender system, using GPT-40 2100 as the base model. We report the results both without (\checkmark) and with (\blacklozenge) parallelization of LLM 2101 calls. For the parallelized version, the maximum number of workers was set to 500. The results in 2102 Table 12 highlight that SimUSER and Rec4Agent achieve similar inference speeds. As expected, 2103 parallelizing the LLM calls significantly reduces the inference time for SimUSER (\$). This demonstrates two key points: 1) user interactions can be scaled to a larger number of users (1,000+) while 2104 maintaining a low simulation time, 2) the present approach can be parallelized \clubsuit to further reduce 2105 the running time.

2106 K THUMBNAIL CAPTION EXAMPLE

We provide below an example of a caption generated by the perception module, given as input a poster of the movie "GoldenEye (1995)":

Output Perception Module

Input: Poster of the movie "GoldenEye (1995)"

Caption: The movie poster for "GoldenEye 007" prominently features James Bond, portrayed by Pierce Brosnan, in the center, aiming a gun directly at the viewer, conveying a sense of intense action and suspense. To the left, a close-up of a mysterious woman with a confident expression adds intrigue, hinting at her pivotal role in the plot. A helicopter and a fighter jet, both mid-flight, further emphasize the film's thrilling, fast-paced narrative. The iconic "007" logo, with the gun symbol integrated into the number, is prominently displayed at the bottom, reinforcing the film's connection to the legendary spy franchise. The overall mood of the poster is intense and action-packed, appealing to fans of the James Bond series and those looking for a thrilling cinematic experience.

2119 2120 2121

2111 2112

2113

2114

2115

2116

2117

2118

2122 2123

L ETHICS AND REPRODUCIBILITY STATEMENT

2124 This paper proposes an LLM-empowered agent framework designed to simulate user interactions 2125 with recommender systems in a realistic and cost-effective manner. While our approach offers sig-2126 nificant benefits in terms of scalability and efficiency, it also raises ethical considerations. The use of 2127 such agents could lead to unintended consequences, such as bias amplification, where the synthetic 2128 agents might inadvertently reinforce existing stereotypes or present skewed recommendations due 2129 to biases in the training data. Additionally, there is a risk of manipulation of user preferences, as the synthetic agents could be used to subtly influence user behavior by consistently promoting cer-2130 tain types of content without explicit user consent. Furthermore, simulating interactions at a broad 2131 scale could result in the identification and exploitation of behavioral patterns that might encourage 2132 specific user behaviors, potentially leading to negative societal impacts. Finally, there is a concern 2133 that developers or designers might use synthetic users and displace the role of humans and system 2134 stakeholders in the design process. 2135

Regarding privacy, SimUSER leverages LLMs' reasoning ability and factual knowledge without
finetuning LLMs, ensuring that LLMs do not retain or remember user-specific data. In addition,
employing matched personas instead of user-specific information available such as the occupation
or gender prevents the leakage of user privacy information through external APIs.

Furthermore, since our experiments and analyses are conducted in English, we do not claim that our
findings are universally applicable across all languages. However, the SimUSER framework may be
adaptable to other languages with suitable modifications.

We acknowledge these risks and advocate for responsible deployment, including transparency, user awareness, and ongoing monitoring to mitigate potential harm. We also suggest that synthetic uses should not be a substitute for real human input in studies and design processes. Rather, these agents should be leveraged during the initial design phases to explore concepts, especially in situations where recruiting human participants is impractical or where testing certain theories with real people could be challenging or pose risks. By adhering to these principles, we can ensure that the deployment of synthetic users in the wild is ethical and socially responsible.

In the interest of reproducibility, we have evaluated SimUSER and baseline methods using publicly available datasets and codebases. The code and scripts necessary to reproduce SimUSER and experiments are available at https://github.com/SimUSER-paper/SimUSER.

- 2153 2154
- 2155
- 2156
- 2157

2158