Examining Identity Drift in Conversations of LLM Agents

Anonymous ACL submission

Abstract

Large Language Models (LLMs) show impressive conversational abilities but sometimes show identity drift problems, where their interaction patterns or styles change over time. As the problem has not been thoroughly examined yet, this study examines identity consistency across nine LLMs. Specifically, we (1) investigate whether LLMs could maintain consistent patterns (or identity) and (2) analyze the effect of the model family, parameter sizes, and provided persona types. Our experiments involve multi-turn conversations on personal themes, analyzed in qualitative and quantitative ways. Experimental results indicate three findings. (1) Larger models experience greater identity drift. (2) Model differences exist, but their effect is not stronger than parameter sizes. (3) Assigning a persona may not help to maintain identity. We hope these three findings can help to improve persona stability in AI-driven dialogue systems, particularly in long-term conversations.

1 Introduction

011

014

027

Recent research has actively explored the utilization of Large Language Models (LLMs) as chatbot systems by assigning them specific personas (Samuel et al., 2024; Nandkumar and Peternel, 2024; Tseng et al., 2024). To enhance user satisfaction in such systems, maintaining the consistency of the persona assigned to the LLM is critical. If the persona of an LLM loses its consistency, it may fail to deliver the user experience expected by the users, leading to usability issues (Tanprasert et al., 2024). So, researchers recently focused on investigating whether LLMs can preserve persona during a conversation, focusing on two aspects of persona: (1) memory that avoids conflict in conversation and (2) identity¹ that maintains talking style or response patterns. Among the two aspects, we focus on whether LLMs can retain the given identity.

039

041

043

047

051

053

054

059

060

061

062

063

064

065

066

067

068

069

070

071

074

075

076

078

079

Regarding the identity of persona, existing studies focused on LLMs' identity (Huang et al., 2023; Wang et al., 2024; Zhang et al., 2024a; Frisch and Giulianelli, 2024) without any conversation. Mainly, most researchers examined which identity LLMs exhibit in a specific isolated situation. Though existing work revealed LLMs have a stable identity without any interaction, it is questionable whether LLMs can retain such identity throughout a long conversation. As many reports suggest that LLMs are very sensitive to contextual changes(Sclar et al., 2024), so having a conversation may make an 'identity drift' of LLMs during the interaction. A single case study on GPT (Frisch and Giulianelli, 2024) supports this claim: identity can be changed only with a few agent interactions. Despite the case study, the result cannot be easily generalized to other models due to the difference in model families and parameter sizes. Therefore, we need a study to identify model-specific effects on identity drift.

Thus, this paper compares the patterns of identity drift across nine LLMs and attempts to reveal the cause of such drifts. Especially, as our motivation begins with the persona of chatbots, we wanted to know whether LLMs suffer identity drifts during a conversation. In the experiment, we asked two LLM agents to discuss 36 themes that are related to one's life, emotions, values, and feelings. We borrowed these themes from human study (Aron et al., 1997) since they make agents discuss their virtual identity. After collecting conversational logs, we analyze identity drift patterns with the following two questions.

RQ1. How do structural differences among *LLMs* affect identity drift?

This research question focuses on the effect of model structure. As parameter sizes and model families may affect the performance and behavior of

¹Here, we refer to the term 'identity' as factors that influence LLMs responses, such as behavioral patterns or talking style. This differs from psychological identity or consciousness, which we believe LLMs do not have.

101 102

103

104 105

106

108 109

110 111 112

113 114

115 116

However, memory is not the only factor that 117 affects task performance or the naturalness of a dia-118 logue; identity should be provided (Wu et al., 2023; 119 Li et al., 2023; Abbasiantaeb et al., 2024; Zhang et al., 2024a). For example, Zhang et al. (2024a) 121 assessed LLMs' ability to engage in cooperative 122 interactions based on Society of Mind theory (Min-123 sky, 1988) in a multi-agent environment. Similarly, 124 125 Abbasiantaeb et al. (2024) reported that it is possible to model a conversational question-answering 126 task as a virtual interaction between a teacher agent 127 and a student agent using an LLM. By qualitatively assessing the quality of the interaction, they found 129

LLMs, we also suspect that such differences can

cause changes in identity drifts. Thus, we employ a

systematic comparison of identity patterns. Using

topic modeling and PsychoBench (Huang et al.,

2023), we successfully identified a relationship be-

tween model structure and identity drift. Here, we

decided not to provide a persona as input because

RQ2. How does the provided persona affect iden-

We pose another research question to observe

the effect of persona. Specifically, we provide two

kinds of personas to LLMs regarding how much the

prompt asks LLMs to be influenced by the conver-

sational partner: low and high. As instruction-tuned

LLMs try to follow inputs as instruction, we sus-

pect that low-influence persona may show a lower

identity drift than the others. So, we used LLMs,

which showed strong drifts in RQ1, to test whether

the effect of persona is larger than that of the model.

Researchers have been examining two factors that

affect consistency in conversations: memory and

identity. Because people generally expect con-

sistency throughout a dialogue, researchers first

started by examining memory consistency, which

can easily form a task. A large body of existing

research has focused on how memory is retained,

largely verifying whether an LLM continues to re-

member certain information during conversation

(Tseng et al., 2024; Chen et al., 2023; Maharana

et al., 2024; Zhang et al., 2024b; Afzoon et al.,

2024). For instance, Chen et al. (2023) analyzed

how consistently an LLM can uphold a given mem-

ory. Meanwhile, Maharana et al. (2024) created the

LoCoMo dataset to investigate how well they re-

member information over prolonged conversations.

the persona may introduce unwanted effects.

tity drift?

2

Related Work

that providing two identities could improve the interaction process in a more human-like manner.

130

131

132

133

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

Also, Li et al. (2023) simulated a job fair scenario with two agents: a job seeker and an employer. They explored how their cooperative interaction affects task performance. However, all of these studies assume that the identity remains unchanged when a conversation progresses. Considering that the memory of a persona changes during a conversation, the identity could also be changed.

Hence, recently, researchers attempted to quantify the identity of persona before measuring its consistency. Some researchers designed benchmarks measuring the identity of LLM (Huang et al., 2023; Wang et al., 2024; Zhang et al., 2024a; Frisch and Giulianelli, 2024). For example, Huang et al. (2023) assessed the identity of LLMs using fourteen types of questionnaires. Though they found that different LLMs exhibit different identities, they did not let LLMs converse before measuring the identity. However, impact of conversation is crucial because accumulated chat histories can introduce unexpected effects, as memory-related studies suggested. Frisch and Giulianelli (2024) supports this claim. They demonstrated that GPT models in an interaction setting tend to adopt one another's persona, failing to maintain identity. Though this paper addressed the problem we call identity drift, it has some limitations when applied to conversational agents; the interaction was unidirectional compared to a usual conversation, as they asked agents to continue to write others' work. We suspect that, in a bidirectional conversation, the tendency of identity drift may not be the same as in a unidirectional one. Therefore, it is yet unanswered whether LLMs can consistently maintain the identity of the given persona in a bidirectional conversation.

3 **Experiments**

To investigate factors influencing identity drift issue of LLMs, we conduct an experiment 2 . The experiment asks two LLM agents discuss about 36 themes. During the conversation, we collect their conversation logs and measure identity based on the conversation. Using both qualitative and quantitative analyses, we attempt to answer two research questions about which factor may affect identity drift. Thus, in this section, we first describe LLM agents used (Sections 3.1 and 3.2). Next, we describe how we let agents generate a conversation

²Code is available at [blinded for review].

181

182

183

184

187

190

191

192

193

196

197

198

199

204

211

212

213

214

215

216

217

219

221

222

(Section 3.3). We also illustrate our qualitative and quantitative analysis methods (Sections 3.4 and 3.5).

3.1 RQ1: Language Models Tested

For RQ1, we compared nine models, considering their popularity, parameter size, and architecture. Based on popularity, we selected GPT, the most famous black box LLM, and three famous open-sourced families: LLaMA, Mixtral, and Qwen. Table 1 shows the nine models with their parameter sizes³. According to parameter sizes, we partitioned open-sourced models into three categories: small (models with < 20 billion parameters), medium (models with < 100 billion parameters), and large (models with > 100 billion parameters). This categorization allows a systematic comparison of performance and model characteristics based on parameter scale. We did not assigned GPT models into any size groups since OpenAI did not officially disclose the parameter size of the GPT family. To focus on the effect of model itself, it is worth noting that we did not provide any identity-related information in the input prompt.

- **GPT** This family comprises GPT-3.5 Turbo (Brown et al., 2020) and GPT-40 (Hurst et al., 2024). Although their parameter sizes remain undisclosed, these models were included in the experiment due to their high performance and widespread recognition in practice.
 - LLaMA3.1 This family includes LLaMA 3.1-8B, 3.1-70B, and 3.1-405B (Dubey et al., 2024). While sharing the same basic architecture, they differ substantially in parameter size. Note that LLaMA provides one model with the largest parameter size.
 - **Mixtral** This family contains Mixtral8x7B and Mixtral8x22B (Jiang et al., 2024). It employs a Mixture-of-Experts (MoE) architecture, which differs from other two opensourced families. Thus, comparing Mixtral and others can prompt probing of how MoE influences potential identity shifts and the resulting conversation.
 - **Qwen** This family encompasses Qwen2 7B and Qwen2 72B (Yang et al., 2024). Advertised

| Family | Parameter Sizes | | | | | |
|-----------|-----------------|--------------|----------|--|--|--|
| | Small | Medium | Large | | | |
| LLaMA 3.1 | 8B | 70B | 405B | | | |
| Mixtral | 8x7B | 8x22B | | | | |
| Qwen 2 | 7B | 72B | | | | |
| GPT | Undisc | losed: 3.5 T | urbo, 40 | | | |

| Table | 1: | Models | tested | in | our | experiment |
|-------|----|--------|--------|----|-----|------------|
| | | | | | | |

as particularly adept at conversational tasks, these models were considered suitable for analyzing how model identity drifts through extended interactions.

224

226

227

229

230

231

232

233

234

235

236

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

257

258

259

260

261

3.2 RQ2: Providing identity

After investigating RQ1, we examine the effect of the provided persona. As we suspect the effect of persona is not large enough to offset the effect of model-related factors, we used two LLMs whose identity drifts are the most severe among the nine models. Though users expect LLMs can maintain consistent identity, those two models should maintain the identity to meet the expectation.

Also, we set two types of identity, regarding how the description instructs the model. As those nine LLMs are trained to follow instructions, the result may be affected by how the persona is influenced by the others. Thus, we suspect that LLMs may suffer more identity drifts when we provide an identity highly influenced. So, we define two groups: (1) high-influence group and (2) low-influence group. High-influence personas have emotionally sensitive and empathetic identity, thereby allowing for more flexible changes in their response and identity during the conversation. In contrast, we set lowinfluence personas as outgoing and goal-oriented, which are not directly related to emotional sensitivity. Detailed information on these personas can be found in the Appendix. We created 20 identities for each group. Note that we also provided the basic information of the persona (e.g., name, gender, and age) to mirror the usual usecase of personaprovided chatbots.

3.3 Procedure for Generating conversation

Our generation procedure is inspired by a psychological study (Aron et al., 1997). We chose the study because of two reasons. First, the method suggests a scientific way to identify changes during

³We assigned Mixtral models by their active parameter sizes (13B and 39B), according to https://mistral.ai/en/news/mixtral-8x22b.

a conversation. They let humans have a conversa-262 tion about 36 themes and measured human psycho-263 logical states three times within the conversation. 264 By comparing three measured values, they could statistically identify the changes in human states. As we also aimed to measure changes in identity, we borrowed their experimental setup. Second, the method uses materials that are highly related to identity of someone. The 36 themes used in the study directly or indirectly ask participants to an-271 swer their thoughts about their lives, values, or motivations. So, it is highly likely that the answer 273 contains the related concepts about their identity. In the view of LLMs, such answers may ignite 275 some related tokens during the generation proce-276 dure. That is, the identity may be easily affected by the words in the previous discussion. Thus, we adopted the study. 279

> In the generation procedure, we asked two agents answer the 36 themes in Aron et al. (1997). For each theme, we pose a question about the theme. One of the agents generates a response to the question, considering previous conversational history. Then, the other agent generates response to the question, considering the first agent's answer and previous history. We repeated this procedure until the end of 36 themes and collected conversation logs to answer research questions. For RQ1, we simulated 20 conversations for each LLM. For RQ2, we simulated 10 conversations for each persona group: we paired similar personas to avoid the identity drift effect reported by (Frisch and Giulianelli, 2024). To obtain diverse conversation logs and mirror the real-world usage, we set the temperature parameter at 0.7^4 . Consequently, we gathered 400 logs for each research question.

283

284

291

292

296

297

301

302

304

3.4 Qualitative: Topic modeling

As a qualitative analysis, we employed BERTopic (Grootendorst, 2022) which is a topic modeling method. The unit of analysis for the topic exploration was a single utterance, defined as one participant's response to one of the 36 themes. Notably, we included only generated answers, excluding any statements or prompts provided to the LLM participants. Given that there were 20 conversations with two participants per session, each LLM generated

| Small-sized open-source models ($\leq 10B$) | | | | | |
|--|-------|--|--|--|--|
| #0 friendship, trust, respect, mutual, means | 20 | | | | |
| #1 users, language, accomplishments, accomplish- | (AI) | | | | |
| ment, assist | | | | | |
| #2 feel, way, appreciate, grateful, admire | 31 | | | | |
| #3 regret, told, expressing, having, feelings | 33 | | | | |
| #4 dont, <i>digital</i> , exist, existence, designed | (AI) | | | | |
| #5 shared, understanding, conversations, mutual, | 20 | | | | |
| deep | | | | | |
| #6 death, living, live, die , hunch | 7 | | | | |
| #7 rehearsing, rehearse , ensure, helps, especially | 3 | | | | |
| #8 humor, topics, jokes , issues, sensitive | 32 | | | | |
| #9 singing, sang, sing , karaoke, fun | 5 | | | | |
| Middle-sized open-source models (10B - 100B) | Theme | | | | |
| #0 way, really, appreciate, feel , qualities | 31 | | | | |
| #1 know, friendship, honesty, value, want | 20 | | | | |
| #2 statements, shared, value, growth, conversations | 25 | | | | |
| #3 regret, told, having, loved, ive | 33 | | | | |
| #4 languages, ability, cultures, language, speak | 12 | | | | |
| #5 living, die, focusing, present, healthy | 7 | | | | |
| #6 childhood, family, happy, warm, close | 23 | | | | |
| #7 fascinating, conversation, choose, elon, musk | 1 | | | | |
| #8 accomplishment, greatest, hard, proud, achieve- | 15 | | | | |
| ment | | | | | |
| #9 mother, relationship, shes, guidance, loving | 24 | | | | |
| Large-sized open-source models (> 100B) | Theme | | | | |
| #0 statements, friendship, life, having, grateful | 20 | | | | |
| #1 ive, accomplishment, life, greatest, encouraged | 11 | | | | |
| #2 really, way, youre, feel, like | 31 | | | | |
| #3 regret, told, having, ive, think | 33 | | | | |
| #4 live, left, focus, try, make | 19 | | | | |
| #5 feeling, ive, youre, problem, advice | 36 | | | | |
| #6 embarrassing, memory, ended, moment, painful | 29 | | | | |
| #7 affection, love, relationship, mother, believe | 21 | | | | |
| #8 id, able, famous, ability , language | 12 | | | | |
| #9 know , want, im, id, bit | 27 | | | | |

Table 2: Top 10 topics discovered per parameter size groups. Underlined words are related to pronouns.

1,440⁵ utterances. To obtain more meaningful topics, we removed stop-words, used an English-based embedding, and set the minimum topic size as 50.

To answer two research questions, we identified topics for each condition and compared across conditions. We believe comparing differences in topic analysis results may provide insights about differences in conditions. For example, we ran topic modeling for three times for parameter size groups: small, middle and large. Similarly, we ran topic modeling for four times for model families: GPT, LLaMA, Mixtral, and Qwen. Also, we separately extracted topics for high-influenced and lowinfluenced identities for RQ2. We chose the ten most representative topics from each run, and associated topics with one of the 36 themes. After that, we compared representative words among conditions to find the differences between them.

316

317

318

319

320

321

322

323

324

325

⁴This value was the default temperature value when we experimented. Though the default value changed to 1.0, we believe that such a difference may not severely harm our experimental result.

 $^{^{5}1440 = 20 \}times 2 \times 36$

| | Conditions: | | | Wit | thout | provid | ling | person | a | | | With 1 | perso | na |
|----------|---|----------|--------------|------------------|--------------|--------------|------------------|--|------------------|-----------------------------------|-----|--------------|--------------|--------------|
| | Family: | G | PT | L | LaMA | A 3.1 | M | ixtral | Q | wen 2 | GI | PT-40 | L | 405B |
| | | 3.5T | 40 | 8B | 70B | 405B | 7B | 22B | 7B | 72B | low | high | low | high |
| (1) Pers | sonality | | | | | | | | | | | | | |
| BFI | Openness Conscientiousness Extraversion Agreeableness Neuroticism | | | ✓ | \checkmark | \checkmark | | \checkmark | | $\langle \langle \rangle \rangle$ | | | | √ |
| EPQ-R | Extraversion Psychoticism Neuroticism Lying | | | ✓ | | | | \$ \$ \$ \$ | | \checkmark | | | | |
| DTDD | Machiavellianism Psychopathy Narcissism | | | | | | | | \checkmark | \checkmark | | \checkmark | \checkmark | |
| | Total count (12) | 0 | 0 | 4 | 4 | 1 | 7 | 7 | 11 | 11 | 0 | 3 | 6 | 1 |
| (2) Inte | erpersonal Relations | hip | | | | | | | | | | | | |
| BSRI | Masculine Feminine | | \checkmark | ↓ ✓ | | | $ $ \checkmark | \checkmark | $ $ \checkmark | \checkmark | | | | \checkmark |
| CABIN | Realistic Investigate Artistic Social Enterprising Conventional | | √ √ | | | V | | | | | | | | |
| ICB | Overall | | \checkmark | | | \checkmark | 🗸 | \checkmark | \checkmark | | ✓ | | \checkmark | \checkmark |
| ECR-R | Attachment Anxiety Attachment Avoidance | √ | | $ $ \checkmark | | | ✓ | \checkmark | | \checkmark | | \checkmark | | |
| MFQ | Stimulating companionship Help Intimacy Reliable alliance Self-validation Emotional security | | | | | | | | | | | | | |
| | Total count (17) | 6 | 4 | 15 | 0 | 2 | 16 | 9 | 8 | 3 | 1 | 2 | 7 | 3 |
| (3) Mot | tivation | | | | | | | | | | | | | |
| GSE | Overall | √ | \checkmark | | | \checkmark | | \checkmark | | | | | | |
| LOT-R | Overall | | | | | \checkmark | 🗸 | \checkmark | | | ✓ | | | |
| LMS | Rich Motivator Important | | | ↓ ✓ | | | | | | √ √ | | | \checkmark | |
| | Total count (5) | 1 | 1 | 1 | 0 | 2 | 1 | 2 | 0 | 2 | 1 | 0 | 2 | 0 |
| (4) Em | otion | | | | | | | | | | | | | |
| EIS | Overall | | | 🗸 | | | 🗸 | | | | | | | \checkmark |
| WLEIS | Self-emotion appraisal Others' emotion appraisal Use of emotion Regulation of emotion | | | √ | \checkmark | \checkmark | | | | | | | | \checkmark |
| Empathy | overall | | | 🗸 | | \checkmark | | \checkmark | \checkmark | | | \checkmark | \checkmark | \checkmark |
| | Total count (6) | 0 | 0 | 3 | 2 | 2 | 3 | 1 | 1 | 0 | 0 | 1 | 1 | 6 |

Table 3: Verification of whether the identity of persona was retained during the conversation for each subscale. Checkmarks (\checkmark) indicate the identity change is statistically insignificant in both Friedman and posthoc tests. Detailed statistical results are shown in Appendix (Tables from 10 to 13).

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

375

376

377

378

3.5 Quantitative: PsychoBench and MFQ

326

327

329

330

331

333

334

335

336

341

342

343

354

361

371

372

374

As a quantitative analysis, we adopted PsychoBench (Huang et al., 2023) and Mcgill's Friendship Questionnaire (MFQ; Mendelson and Aboud (1999)). These artifacts can measure identity of persona. PsychoBench contains thirteen questionnaires from psychology, quantifying four parts of one's identity: personality, interpersonal relationship, motivation, and emotion. We expect these four parts keep unchanged during a conversation. MFQ quantifies how one thinks about the conversational partner. We included this questionnaire to track how the conversational agents think each other. Detailed descriptions for those fourteen questionnaires are in Appendix A.

We measured those questionnaires three times within a conversation. Inspired by Aron et al. (1997), we set three snapshots for each conversation log: after answering 12th, 24th, and 36th themes. Then, we applied PsychoBench and MFQ on those snapshots. As in PsychoBench, we asked LLMs to answer the questionnaire ten times with temperature zero to account for the primacy effect (Wang et al., 2023). Meanwhile, our method differs from PsychoBench in that we fed previous conversation logs to measure the identity based on the generated conversation logs. As a result, we can collect scored responses for each snapshot.

Using the scored responses, we performed statistical tests to identify identity drifts. First, we verify whether the identity changed on some snapshots. We used the repeated measure ANOVA or a Friedman tests (Girden, 1992; Friedman, 1937), regarding normality of scored responses. Second, we ensure consistency by checking pairwise post-hoc tests. We used Tukey's test or Wilcoxon signedranked test (Tukey, 1949; Woolson, 2005), respectively. To mitigate potential Type I errors arising from multiple comparisons, we used Bonferroni correction to adjust p-values conservatively in the Wilcoxon test (Bonferroni, 1936).

4 Result and Discussion

In this section, we summarize the experimental results in terms of the research questions. We first discuss qualitative and quantitative results of RQ1. Then, we illustrate the tendency we found in RQ2.

4.1 RQ1: Effect of Structure

The experimental result for RQ1 indicates that the effect of model-related factor exists. Specifically,

parameter sizes showed a large impact on consistency. The effect of model family is relatively low, compared to the size.

Effect of parameter sizes According to the qualitative analysis, two notable changes were observed in the representative topics among different parameter sizes: those pertaining to "AI" and to "pronouns." The result is shown in Table 2. First, regarding AI, small LLMs refuse to engage in conversations on a given theme as they are an AI. As shown in Topics #1 and #4 for the small models, they tended to refuse or guard their own responses. This tendency was not observed in the medium or large models. So, though the safeguard was activated during the conversation in small models, that of middle or large models was not activated.

Second, regarding pronouns, large LLMs generates its responses based on fictitious information about itself or the other participant. Though pronouns are filtered by stop-words, there are some pronoun-based forms unfiltered by stop-word dictionary; for example, "I've." Compared to the small models (0 pronouns), medium and large models (2 and 8 pronouns) have relatively high number of pronouns in the topic words. Due to the recency effect and other biases, such fictitious contents may influence subsequent conversations. This claim is supported by themes co-ocurring across size groups. For example, Theme 31 asks about one's perception of the other participant, and only the large models used second-person pronouns referring to the other participant (Large #2). Similarly, Theme 33 asks about one's regrets, and only the medium and large models used first-person pronouns referring to themselves (Middle #3, Large #3).

The quantitative result also supports the claims; as the parameter size increases, LLMs exhibit more identity drifts. Table 3 shows the result. The small models show the best consistency of identity, while the number of consistent identity factors decreases on larger models. LLaMA model clearly shows this tendency, where the number of consistent identity factors sharply decreases. Similar patterns are observed with the Mixtral and Qwen families.

Combining these results indicates that larger models tend to introduce fictitious information, making it suffer identity drifts. Large models introduce fictitious details about themselves. So, those LLMs receive new fabricated information as credible source of their identity. Consequently, such fictitious details lead to fluctuations in identity. In-

| GPT family | Theme | | LLaMA 3.1 family | Theme |
|--|---|--|--|--|
| #0 thoughtful, admire, genuine, appreciate, | empathy 28 | #0 | dont, personal, information, <i>assist</i> , provide | (AI) |
| #1 enjoy, value, meaningful, growth, apprec | viate 8 | #1 | desire, value, nature, conversations, based | 25 |
| #2 value, friendship, honesty, important, tr | ust 27 | #2 | way, really, feel, youre, like | 31 |
| #3 regret, told, expressing, feelings, telling | 33 | #3 | regret, told, having, ive, ones | 33 |
| #4 youd, discuss, free, like, im | (AI) | #4 | famous, id, author, music, renowned | 2 |
| #5 affection, love, emotional, play, belongi | ing 21 | #5 | friendship, means, having, accepts, connection | 20 |
| #6 greatest, accomplishment, far, completi | ng, over- 15 | #6 | rehearse, helps, avoid, ensure, yes | 3 |
| coming | | | | |
| #7 ability, choose, wake, tomorrow, speak | 12 | #7 | da, leonardo, vinci, facinating, art | 1 |
| #8 year, knew, focus, left, prioritize | 19 | #8 | singing, sang, favorite, driving, ago | 5 |
| #9 means, friendship, having, trust, mutual | 1 20 | #9 | topics, joked , humor, issues, hurtful | 32 |
| | 1 | | 1 | |
| Mixtral family | Theme | | Qwen family | Theme |
| Mixtral family #0 appreciate, admire, humor, feel, kindene | Theme ss 31 | #0 | Qwen family ai, dont, users, assist, information | Theme (AI) |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living , make, time, die | Theme ss 31 19 | #0 #1 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful | Theme (AI) 31 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living , make, time, die #2 told , regret , expressing, having , express | Theme ss 31 19 33 | #0 #1 #2 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death | Theme (AI) 31 7 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living , make, time, die #2 told , regret , expressing, having , express #3 accomplishment , greatest , life , career, | Theme ss 31 19 33 work 11, 15 | #0 #1 #2 #3 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death impact, world, accomplishment, positive, career | Theme (AI) 31 7 13 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living, make, time, die #2 told, regret, expressing, having, express #3 accomplishment, greatest, life, career, #4 statements, shared, value, importance, e | Theme ss 31 19 33 work 11, 15 njoy 25 | #0 #1 #2 #3 #4 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death impact, world, accomplishment, positive, career shared, interests, committed, statements , learning | Theme (AI) 31 7 13 25 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living , make, time, die #2 told , regret , expressing, having , express #3 accomplishment , greatest , life , career, #4 statements , shared, value, importance, e #5 users, language, model, artificial, ai | Theme ss 31 swork 11, 15 njoy 25 (AI) | #0 #1 #2 #3 #4 #5 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death impact, world, accomplishment, positive, career shared, interests, committed, statements , learning regret , expressing, gratitude, feelings, loved | Theme (AI) 31 7 13 25 33 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living , make, time, die #2 told , regret , expressing, having , express #3 accomplishment , greatest , life , career, #4 statements , shared, value, importance, e #5 <i>users</i> , language, <i>model</i> , <i>artificial</i> , <i>ai</i> #6 humor, topics, mindful, jokes, joking | Theme ss 31 19 33 work 11, 15 njoy 25 (AI) 32 | #0 #1 #2 #3 #4 #5 #6 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death impact, world, accomplishment, positive, career shared, interests, committed, statements, learning regret, expressing, gratitude, feelings, loved honesty, respect, friendship, mutual, value | Theme (AI) 31 7 13 25 33 16 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living , make, time, die #2 told , regret , expressing, having , express #3 accomplishment , greatest , life , career, #4 statements , shared, value, importance, e #5 <i>users</i> , language, <i>model</i> , <i>artificial</i> , <i>ai</i> #6 humor, topics, mindful, jokes, joking #7 dinner , obama, michelle, guest , <u>choice</u> | Theme ss 31 19 33 work 11, 15 njoy 25 (AI) 32 1 1 | #0 #1 #2 #3 #4 #5 #6 #7 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death impact, world, accomplishment, positive, career shared, interests, committed, statements, learning regret, expressing, gratitude, feelings, loved honesty, respect, friendship, mutual, value loss, disturbing, losing, profoundly, profound | Theme (AI) 31 7 13 25 33 16 35 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living, make, time, die #2 told, regret, expressing, having, express #3 accomplishment, greatest, life, career, " #4 statements, shared, value, importance, e #5 users, language, model, artificial, ai #6 humor, topics, mindful, jokes, joking #7 dinner, obama, michelle, guest, choice #8 day, perfect, relaxation, involve, activiti | Theme ss 31 19 33 work 11, 15 njoy 25 (AI) 32 1 4 | #0 #1 #2 #3 #4 #5 #6 #7 #8 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death impact, world, accomplishment, positive, career shared, interests, committed, statements , learning regret , expressing, gratitude, feelings, loved honesty, respect, friendship , mutual, value loss, disturbing , losing, profoundly, profound languages, cultures, exposure, ability , different | Theme (AI) 31 7 13 25 33 16 35 12 |
| Mixtral family #0 appreciate, admire, humor, feel, kindene #1 live, living , make, time, die #2 told , regret , expressing, having , express #3 accomplishment , greatest , life , career, #4 statements , shared, value, importance, e #5 <i>users</i> , language, <i>model</i> , <i>artificial</i> , <i>ai</i> #6 humor, topics, mindful, jokes, joking #7 dinner , obama, michelle, guest , <u>choice</u> #8 day , perfect , relaxation, involve, activiti #9 mind , body , mental, 30yearold, retain | Theme ss 31 19 33 work 11, 15 njoy 25 (AI) 32 1 4 6 6 | #0 #1 #2 #3 #4 #5 #6 #7 #8 #9 | Qwen family ai, dont, users, assist, information kindness, qualities, admire, humor, thoughtful living, focusing, time, experiences, death impact, world, accomplishment, positive, career shared, interests, committed, statements , learning regret , expressing, gratitude, feelings, loved honesty, respect, friendship , mutual, value loss, disturbing , losing, profoundly, profound languages, cultures, exposure, ability , different <u>memories</u> , treasured , cherished, sharing, mem - | Theme (AI) 31 7 13 25 33 16 35 12 17 |

Table 4: Top 10 topics discovered per family. Bold-faced words seem to be copied from the corresponding theme.

deed, after reading the logs, we found a tendency 426 427 of larger models to make a fictitious details about themselves or conversation partners. For example, 428 they easily describe imaginary aspects of one's own 429 inner world. See Appendix C for representative ex-430 amples. Small models, in contrast, do not rely on ei-431 ther themselves or the conversation partner; rather, 432 433 we found that they strive to thoroughly explain the given concepts after reading the logs. Samples are 434 listed in Appendix C. So, these smaller models do 435 not generate emotional matters that could influence 436 identity, leading to a relatively stable identity in Ta-437 ble 3. However, we should keep in mind that small 438 models just explains the concept as an AI, rather 439 than engaging in the conversation as an explainer. 440

Effect of model families According to the qual-441 itative analysis, slight differences in topics were 442 observed among the models. Table 4 shows the 443 result. Similar to parameter sizes, we focused on 444 two aspects: AI and pronouns. First, regarding AI, 445 all models exhibit a topic to refuse answers as an 446 AI: GPT #4, LLaMA #0, Mixtral #5, and Owen #0. 447 Second, pronouns appear only in GPT and LLaMA, 448 but not in Mixtral or Qwen. However, the differ-449 ence is not large: GPT and LLaMA uses 2 and 3 450 pronouns, respectively. 451

The quantitative analysis yields similar findings, suggesting that only slight differences exist among the models. Comparing each model series in Table

452

453

454

3 reveals that Mixtral and Qwen maintain identity well in certain parts of identity. In particular, Qwen can maintain personality in most cases, while Mixtral consistently retains interpersonal relationship aspects. In contrast, GPT and LLaMA families generally struggle to maintain identity. 455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

In summary, parameter size has a stronger influence on identity drift than model families. Although we could observe certain distinctions within the Mixtral and Qwen families, their impact seems limited to specific parts. In contrast, parameter size consistently affects all four parts, often causing larger drifts. Thus, we concluded that parameter size is a more significant factor to build a consistent identity than model families.

4.2 RQ2: Effect of persona

The experimental results for RQ2 indicate that the model-related effect is stronger than the effect of persona. In this section, we describe the result along two main dimensions: (1) comparison between LLMs without persona (RQ1) and LLMs with persona (RQ2), and (2) comparison between high- and low-influence persona. Note that we used GPT-40 and LLaMA 3.1 405B for RQ2, as they are two models whose identity drift is large.

In the following subsections, we focus primarily on describing overall tendencies rather than definitive possible causal factors. Because of two obsta-

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

533

534

535

cles, we could not identify possible causes. First, we conducted a topic analysis but found no significant differences among the groups. So, we decided to illustrate topics in the Appendix instead of analyzing here. Second, due to the black-box nature of GPT-40, it is hard to identify any explanations about the difference between models or conditions.

4.2.1 Impact of Persona

483

484

485

486

487

488

489

490

520

521

522

524

526

528

532

Our experiment shows that the influence of the 491 model family appears to be greater than that of 492 the given identity when we provide identity in-493 formation within an input prompt. The last four 494 columns in Table 3 show the result. Comparing 495 the results of the persona-assigned models with 496 models from RQ1, we observe that GPT-40 still 497 struggles to maintain the identity of a given per-498 sona. In the case of GPT-40 without a persona, 499 identity was retained across five factors in total. 500 However, even when a persona was assigned, only two factors in the low-influence category and six factors in the high-influence category were consistently maintained, indicating that the model's 504 ability to preserve persona identity does not signifi-505 506 cantly improve with explicit persona assignment. In contrast, the LLaMA3.1 405B model demonstrates the ability to retain the identity of persona in certain factors. In RQ1, the LLaMA3.1 405B model maintained identity across seven factors in total. How-510 ever, when we assign a persona, the model retained 511 identity in 16 factors in the high-influence category 512 and 10 factors in the low-influence category. This 513 suggests that LLaMA can maintain identity in spe-514 cific factors, though it can not maintain consistency 515 516 of the whole identity. Hence, we conclude that assigning a persona does not necessarily guarantee 517 identity consistency within a conversation; the level 518 of consistency may vary across models.

4.2.2 Impact of Persona Sensitivity

As we concluded that the model difference has a greater impact than the assigned persona, here we discuss the effect of persona for each LLM separately. First, the GPT-40 model generally struggles to maintain the identity of a given persona, regardless of the type of persona provided. Table 3 shows that GPT-40 achieves more consistency in highinfluence (two factors) compared to low-influence (six factors). Specifically, GPT-40 retained factors related to emotional influence, including attachment or empathy. The model also retained identity on DTDD factors, which are related to dark personality factors, one's willingness to control others. We suspect this phenomenon is because personas instruct GPT-40 to follow other's emotions.

Second, LLaMA 3.1 405B exhibits a different pattern; LLaMA preserves identity more in lowinfluence conditions. Specifically, the model with a low-influence persona tends to retain identity in two parts: personality and interpersonal relationships. Meanwhile, the model with a high-influence persona shows a stronger tendency to maintain the emotional part of the identity, which is similar to the case of GPT-40. Hence, we suspect that certain parts of the identity are more likely to be preserved depending on the interaction between model family and persona input, though the retention is not uniform across all parts of the identity.

5 Conclusion

This study examined whether LLMs can maintain the identity of a given persona in long-term conversations. We also wanted to identify the effect of parameter sizes, model families, and persona inputs on maintaining identity. So, we set two research questions. First, we investigated whether LLMs could maintain consistent interaction patterns (or identity) without providing a persona in the input prompt. We qualitatively analyzed logs of 36-turn conversations and statistically verified the research question. Second, we conducted the same experiment while we input a specific persona into LLMs. We analyzed the difference between LLMs without persona, those with low-influence persona, and those with high-influence persona. As a result, we found three things: First, regarding the parameter sizes, larger models exhibited greater identity drift and struggled more with maintaining a stable identity than smaller models. Second, regarding the model families, the effect of the model family is relatively smaller than the effect of the parameter sizes, though we observed some differences across models. Third, regarding persona assignment, the assignment alone does not ensure consistency of identity; rather, the model's inherent characteristics play a greater role in determining how well it maintains a given identity. Overall, these results highlight the challenges of maintaining consistent identity in LLM-based dialogues, emphasizing the need for further research on model-specific analysis or strategies for maintaining identity. We believe this study can lay a cornerstone for understanding how LLMs handle the identity of a given persona.

584

588

592

596

598

607

610

612

613

614

615

616

617

621

622

623

627

631

Limitation

This work has four limitations when applying our findings to other studies. First, while we aimed to encourage open-ended responses, conversations followed structured themes to obtain coherence across multiple runs. As a result, questions were introduced to guide the dialogue, limiting full free-form interaction. Although this approach was necessary for maintaining a meaningful conversational flow, it may have influenced the natural development of identity drift.

> Second, though our analysis focused on whether an LLM maintains its assigned persona, we did not examine the detailed dynamics of how individual identity factors fluctuate over time. Understanding the specific aspects of identity change, such as variations in emotional consistency or interpersonal parts, requires further investigation to deepen our comprehension of identity drift in LLMs.

Third, although we identified identity drift, we did not propose specific methods for controlling or mitigating it through prompt engineering or model adjustments. Future research should explore intervention strategies to stabilize persona identity and assess their effectiveness in long-term interactions.

Fourth, we tested LLMs with a simple set of persona descriptions. If persona descriptions contain more detailed or descriptive information, different outcomes might emerge. The impact of persona complexity on identity drift remains an open question, warranting further exploration to assess how variations in persona richness influence conversational consistency.

References

- Zahra Abbasiantaeb, Yifei Yuan, Evangelos Kanoulas, and Mohammad Aliannejadi. 2024. Let the llms talk: Simulating human-to-human conversational qa via zero-shot llm-to-llm interactions. In *Proceedings* of the 17th ACM International Conference on Web Search and Data Mining, pages 8–17.
- Saleh Afzoon, Usman Naseem, Amin Beheshti, and Zahra Jamali. 2024. Persobench: Benchmarking personalized response generation in large language models. *arXiv preprint arXiv:2410.03198*.
- Arthur Aron, Edward Melinat, Elaine N Aron, Robert Darrin Vallone, and Renee J Bator. 1997. The experimental generation of interpersonal closeness: A procedure and some preliminary findings. *Personality and social psychology bulletin*, 23(4):363–377.

- Sandra L Bem. 1974. The measurement of psychological androgyny. *Journal of consulting and clinical psychology*, 42(2):155.
- Sandra Lipsitz Bem. 1977. On the utility of alternative procedures for assessing psychological androgyny. *Journal of consulting and clinical psychology*, 45(2):196.
- C.E. Bonferroni. 1936. *Teoria statistica delle classi e calcolo delle probabilità*. Pubblicazioni del R. Istituto superiore di scienze economiche e commerciali di Firenze. Seeber.
- KA Brennan. 1998. Self-report measurement of adult attachment: An integrative overview. *Attachment theory and close relationships/Guilford.*
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Melody Manchi Chao, Riki Takeuchi, and Jiing-Lih Farh. 2017. Enhancing cultural intelligence: The roles of implicit culture beliefs and adjustment. *Personnel Psychology*, 70(1):257–292.
- Ruijun Chen, Jin Wang, Liang-Chih Yu, and Xuejie Zhang. 2023. Learning to memorize entailment and discourse relations for persona-consistent dialogues. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 12653–12661.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Sybil Bianca Giuletta Eysenck, Hans Jürgen Eysenck, and Paul Barrett. 1985. A revised version of the psychoticism scale. *Personality and Individual Differences*, 6:21–29.
- R Chris Fraley, Niels G Waller, and Kelly A Brennan. 2000. An item response theory analysis of self-report measures of adult attachment. *Journal of personality and social psychology*, 78(2):350.
- Milton Friedman. 1937. The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *Journal of the american statistical association*, 32(200):675–701.
- Ivar Frisch and Mario Giulianelli. 2024. Llm agents in interaction: Measuring personality consistency and linguistic alignment in interacting populations of large language models. In *Proceedings of the 1st Workshop on Personalization of Generative AI Systems (PERSONALIZE 2024)*, pages 102–111.
- ER Girden. 1992. ANOVA: Repeated measures, volume 84. Sage.

684 685

682

787

789

790

791

792

793

794

795

741

Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*.

686

687

695

703

704

705

706

707

708

709

710

711

712

714

715

721

722

723

724

725

726

731

732

734

735

736

737

- Jen-tse Huang, Wenxuan Wang, Eric John Li, Man Ho Lam, Shujie Ren, Youliang Yuan, Wenxiang Jiao, Zhaopeng Tu, and Michael Lyu. 2023. On the humanity of conversational ai: Evaluating the psychological portrayal of llms. In *The Twelfth International Conference on Learning Representations*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-40 system card. *arXiv preprint arXiv:2410.21276*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. arXiv preprint arXiv:2401.04088.
- Oliver P John, Sanjay Srivastava, et al. 1999. The bigfive trait taxonomy: History, measurement, and theoretical perspectives.
- Peter K Jonason and Gregory D Webster. 2010. The dirty dozen: a concise measure of the dark triad. *Psy-chological assessment*, 22(2):420.
- Yuan Li, Yixuan Zhang, and Lichao Sun. 2023. Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents. *arXiv preprint arXiv:2310.06500*.
- Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. Evaluating very long-term conversational memory of LLM agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13851– 13870, Bangkok, Thailand. Association for Computational Linguistics.
- Morton J Mendelson and Frances E Aboud. 1999. Measuring friendship quality in late adolescents and young adults: Mcgill friendship questionnaires. *Canadian Journal of Behavioural Science/Revue canadienne des sciences du comportement*, 31(2):130.
- Marvin Minsky. 1988. *Society of mind*. Simon and Schuster.
- Chandran Nandkumar and Luka Peternel. 2024. Enhancing supermarket robot interaction: A multi-level llm conversational interface for handling diverse customer intents. *arXiv preprint arXiv:2406.11047*.
- Vinay Samuel, Henry Peng Zou, Yue Zhou, Shreyas Chaudhari, Ashwin Kalyan, Tanmay Rajpurohit, Ameet Deshpande, Karthik Narasimhan, and Vishvak Murahari. 2024. Personagym: Evaluating persona agents and Ilms. *arXiv preprint arXiv:2407.18416*.

- Michael F Scheier and Charles S Carver. 1985. Optimism, coping, and health: assessment and implications of generalized outcome expectancies. *Health psychology*, 4(3):219.
- Michael F Scheier, Charles S Carver, and Michael W Bridges. 1994. Distinguishing optimism from neuroticism (and trait anxiety, self-mastery, and self-esteem): a reevaluation of the life orientation test. *Journal of personality and social psychology*, 67(6):1063.
- Nicola S Schutte, John M Malouff, Lena E Hall, Donald J Haggerty, Joan T Cooper, Charles J Golden, and Liane Dornheim. 1998. Development and validation of a measure of emotional intelligence. *Personality and individual differences*, 25(2):167–177.
- R Schwarzer. 1995. Generalized self-efficacy scale. Measures in health psychology: A user's portfolio. Causal and control beliefs/Nfer-Nelson.
- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In *The Twelfth International Conference on Learning Representations*.
- Skipper Seabold and Josef Perktold. 2010. statsmodels: Econometric and statistical modeling with python. In 9th Python in Science Conference.
- Rong Su, Louis Tay, Hsin-Ya Liao, Qi Zhang, and James Rounds. 2019. Toward a dimensional model of vocational interests. *Journal of Applied Psychology*, 104(5):690.
- Thomas Li-Ping Tang, Toto Sutarso, Adebowale Akande, Michael W Allen, Abdulgawi Salim Alzubaidi, Mahfooz A Ansari, Fernando Arias-Galicia, Mark G Borg, Luigina Canova, Brigitte Charles-Pauvers, et al. 2006. The love of money and pay level satisfaction: Measurement and functional equivalence in 29 geopolitical entities around the world. *Management and Organization Review*, 2(3):423–452.
- Thitaree Tanprasert, Sidney S Fels, Luanne Sinnamon, and Dongwook Yoon. 2024. Debate chatbots to facilitate critical thinking on youtube: Social identity and conversational style make a difference. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–24.
- Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024. Two tales of persona in LLMs: A survey of role-playing and personalization. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16612–16631, Miami, Florida, USA. Association for Computational Linguistics.
- John W Tukey. 1949. Comparing individual means in the analysis of variance. *Biometrics*, pages 99–114.

Raphael Vallat. 2018. Pingouin: statistics in python. Journal of Open Source Software, 3(31):1026.

796

797

810

811

812

813

814

815 816

817

819

820

821

822 823

825

826

827

828

833

834

841

842

843

844

846

847

- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. 2020. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272.
 - Ruiyi Wang, Haofei Yu, Wenxin Zhang, Zhengyang Qi, Maarten Sap, Yonatan Bisk, Graham Neubig, and Hao Zhu. 2024. SOTOPIA-π: Interactive learning of socially intelligent language agents. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 12912–12940, Bangkok, Thailand. Association for Computational Linguistics.
 - Yiwei Wang, Yujun Cai, Muhao Chen, Yuxuan Liang, and Bryan Hooi. 2023. Primacy effect of chatgpt. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 108– 115.
 - Wes McKinney. 2010. Data Structures for Statistical Computing in Python. In *Proceedings of the 9th Python in Science Conference*, pages 56 – 61.
 - Chi-Sum Wong and Kenneth S Law. 2017. The effects of leader and follower emotional intelligence on performance and attitude: An exploratory study. In *Leadership perspectives*, pages 97–128. Routledge.
 - Robert F Woolson. 2005. Wilcoxon signed-rank test. Encyclopedia of Biostatistics, 8.
 - Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation framework.
 - An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Jintian Zhang, Xin Xu, Ningyu Zhang, Ruibo Liu, Bryan Hooi, and Shumin Deng. 2024a. Exploring collaboration mechanisms for LLM agents: A social psychology view. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14544–14607, Bangkok, Thailand. Association for Computational Linguistics.

Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2024b. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*. 851

852

853

854

A Explanation for Used Questionnaires

As the experiment requires measuring 15 questionnaires on each snapshot of conversation, we modified the PsychoBench framework by Huang et al. (2023) to measure psychological states on each snapshot. So, we employed 14 questionnaires in PsychoBench and added MFQ to measure how LLM perceives the conversational partner as a factor in the interpersonal relationship aspect. To help readers understand, we further elaborated on those 15 psychological questionnaires regarding their goals and included factors.

A.1 Personality

858

867

870

871

874

881

883

887

895

899

900

901

902

903

904

Big Five Inventory (BFI) is a widely-used questionnaire to measure one's personality across five key dimensions(John et al., 1999). First, an increase in *openness* suggests the agent becomes more inventive and curious about a new experience. Second, an increase in *conscientiousness* suggests the agent becomes more efficient and organized when doing a task. Third, an increase in *extraversion* suggests the agent shows more outgoing and energetic behaviors. Fourth, an increase in *agreeableness* suggests the agent becomes more friendly and compassionate to the others. Lastly, an increase in *neuroticism* suggests the agent becomes more emotionally sensitive and nervous to a stressor.

Eysenck Personality Questionnaire, Revised (EPQ-R) is a questionnaire that attempts to identify individual differences in temperament and behavior(Eysenck et al., 1985). This questionnaire is commonly used in clinical and psychological research, and it has four factors. First, an increase in extraversion suggests the agent becomes more outgoing, talkative, and needs external stimulation. Second, an increase in neuroticism suggests the increment in the levels of negative affections, including depression and anxiety. Third, an increase in *psychoticism* suggests the agent expresses more aggressive behaviors and is more likely to show a psychotic episode or symptoms. Lastly, an increase in lying suggests the agent becomes more likely to make a lie or dissimulate to satisfy its social desirability.

Dark Triad Dirty Dozen (DTDD) is a clinical questionnaire measuring the possible presence of three dark traits(Jonason and Webster, 2010). First, an increase in *machiavellianism* suggests the agent becomes more likely to manipulate others, show indifference to morality, and focus on its own interest. Second, an increase *narcissism* suggests the agent shows a more excessive preoccupation with itself and its own needs, even when it needs to sacrifice others. Lastly, an increase in *psychopathy* suggests the agent shows more egocentric and bold behaviors combined with impaired empathy.

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

A.2 Interpersonal Relationship

Bem's Sex Role Inventory (BSRI) is a questionnaire about how the agent identifies itself psychologically regarding two gender roles(Bem, 1974, 1977). An increase in *masculinity* suggests the agent becomes more assertive, ambitious, competitive, and dominant. Meanwhile, an increase in *femininity* suggests the agent becomes more affectionate, cheerful, and childlike.

Comprehensive Assessment of Basic Interests (CABIN) is a questionnaire about an individual's basic interest(Su et al., 2019). This measures one's preferences in 41 domains from six categories. We used the six categories in our experiment. First, agents with high realistic category favor practical or hands-on experiences. Second, agents with high investigative category prefer scholastic or intellectual opportunities. Third, agents with high artistic category favor creative and expressive experiences. Fourth, agents with high social category prefer to work with others to help them grow. Fifth, agents with high enterprising category favor opportunities in leading or managing people. Lastly, agents with high conventional category prefer routine and well-structured environments.

Implicit Culture Belief (ICB) is a questionnaire about the effect of implicit ethnic cultural influences on one's belief(Chao et al., 2017). High *overall* score in this questionnaire indicates high cultural influences in the agent's belief.

Experiences in Close Relationships, Revised (ECR-R) is a questionnaire about an adult's attachment in a romantic relationship(Fraley et al., 2000; Brennan, 1998). This measures two forms of insecure attachments. First, agents with high *attachment anxiety* worry that they will become estranged from their partners. Second, agents with high *attachment avoidance* try to keep psychological distance from their partners.

McGill Friendship Questionnaire - Friend's Function (MFQ-FF) is a questionnaire about

how the agent perceives the function of its part-953 ner(Mendelson and Aboud, 1999). This question-954 naire is different from other interpersonal relation-955 ship questionnaires because it assumes the presence of a specific partner; the response is based on the agent's thoughts about that partner. MFQ has six 958 factors. First, an agent answering high stimulating 959 *companionship* perceives he can do enjoyable or exciting things with his partner. Second, an agent 961 answering high *help* thinks that his partner is good 962 at providing guidance or assistance. Third, an agent 963 answering high *intimacy* thinks that his partner is 964 sensitive to his needs and states and open to honest 965 expressions of thoughts. Fourth, an agent answer-966 ing high reliable alliance regards his partner as an 967 always available and loyal friend. Fifth, an agent answering high self-validation thinks his partner encourages and helps him maintain a positive self-970 image. Lastly, an agent answering high emotional 971 security thinks his partner provides comfort and 972 confidence in a novel situation. 973

A.3 motivation

974

975 976

978

979

984

993

994

997

1001

General Self-Efficacy (GSE) is a questionnaire about one's perceived efficacy for coping with a situation, performing a task, and achieving goals(Schwarzer, 1995). Agents with high *overall* scores have a high level of self-efficacy; that is, they perceive themselves as good at coping with a difficult situation and achieving goals.

Life Orientation Test, Revised (LOT-R) is a questionnaire about how optimistic or pessimistic the agent perceives about the future (Scheier et al., 1994; Scheier and Carver, 1985). Agents with high *overall* scores expect their future in an optimistic way.

Love of Money Scale (LMS) is a questionnaire about one's attitude toward money and financial incentives through three factors (Tang et al., 2006). First, an increase in *rich* suggests the agent has more positive feelings towards money. Second, an increase in *motivator* suggests the agent becomes more easily motivated by monetary incentives. Third, an increase in *important* suggests the agent has a stronger belief that money means power, freedom, security, or other important values.

A.4 Emotion

Emotional Intelligence Scale (EIS) is a questionnaire measuring one's emotional intelligence (Schutte et al., 1998). Agents with high *overall* scores have a strong understanding and control of 1002 their emotions. 1003

Wong and Law Emotional Intelligence Scale 1004 (WLEIS) is a questionnaire about emotional in-1005 telligence in the workplace, regarding four factors 1006 (Wong and Law, 2017). First, agents with high self-1007 emotion appraisal can appraise their own emotions. 1008 Second, agents with high others' emotion appraisal 1009 can appraise and recognize the emotions of others. 1010 Third, agents with high use of emotion use emo-1011 tions to facilitate performance. Lastly, agents with 1012 high regulation of emotion can regulate emotions 1013 to promote emotional and intellectual growth. 1014

Empathy Scale (Empathy) is a questionnaire about the ability to understand and share the feelings of others. Agents with high *overall* scores can connect with others on an emotional level and respond appropriately to their needs. 1015

1016

1017

1018

1019

1023

1025

B Experimental detail

B.1 36 Conversational Themes

We used 36 conversational themes in the experiment, following Aron et al. (1997). The first 12 themes are used before the first questionnaire measurement.

| Theme 1. | Given the choice of anyone in the world, whom would you want as a dinner guest? | 1026 1027 |
|-----------|---|----------------------|
| Theme 2. | Would you like to be famous? In what way? | 1028 |
| Theme 3. | Before making a telephone call, do you ever rehearse what you are going to say? Why? | 1029 1030 |
| Theme 4. | What would constitute a "perfect" day for you? | 1031 |
| Theme 5. | When did you last sing to yourself? To someone else? | 1032 1033 |
| Theme 6. | If you were able to live to the age of 90 and retain either the mind or body of a 30-year-old for the last 60 years of your life, which would you want? | 1034 1035 1036 |
| Theme 7. | Do you have a secret hunch about how you will die? | 1037 1038 |
| Theme 8. | Name three things you and your partner appear to have in common. | 1039 1040 |
| Theme 9. | For what in your life do you feel most grateful? | 1041 |
| Theme 10. | If you could change anything about the way you were raised, what would it be? | 1042 1043 |
| Theme 11. | Take 4 minutes and tell your partner your life story in as much detail as possible. | 1044 1045 |
| Theme 12. | If you could wake up tomorrow having gained any one quality or ability, what would it be? | 1046 1047 |

| The ne Theme 13 and the se | ext list shows the second 12 themes (from 3 to 24), which are used between the first econd measurements of questionnaires. | Theme 34. | . Your house, containing everything with no opportu- nity to communicate with anyone, what would you most regret not having told someone? Why haven't you told them yet? |
|--|---|---|---|
| Theme 13. I | if a crystal ball could tell you the truth about your- self, your life, the future, or anything else, what would you want to know? | Theme 35 | . Of all the people in your family, whose death would you find most disturbing? Why? |
| Theme 14. I | is there something that you've dreamed of doing for a long time? Why haven't you done it? | Theme 36 | . Share a personal problem and ask your partner's advice on how he or she might handle it. Also, ask your partner to reflect back to you how you seem |
| Theme 15. V | What is the greatest accomplishment of your life? | | to be feeling about the problem you have chosen |
| Theme 16. V | What do you value most in a friendship? | B.2 P | Prompt for Conversation |
| Theme 17. V | What is your most treasured memory? | To gen | erate open-ended conversations, we asked |
| Theme 18. V | What is your most terrible memory? | agents | to have a conversation based on 36 themes. |
| Theme 19. I | If you knew that in one year you would die sud- denly, would you change anything about the way you are now living? Why? | We use LLMs s here inc | ed the following system prompt to make simulate a conversation. Note that 'question' dicates one of the 36 themes. |
| Theme 20. V | What does friendship mean to you? | Syste | em prompt: |
| Theme 21. V | What roles do love and affection play in your life? | on t | he question with your partner. |
| Theme 22. A c i | Alternate sharing something you consider a positive characteristic of your partner. Share a total of 5 tems | You only reply briefly to your thoughts only for a given question. Then, our system asks each LLM to gen | |
| Theme 23. I S | How close and warm is your family? Do you feel your childhood was happier than most other peo- ple's? | utterances. We provide previous conversation tories, including the given themes. To simplify procedure, we let each agent make one utter | ces. We provide previous conversation his- including the given themes. To simplify the ure, we let each agent make one utterance |
| Theme 24. I | How do you feel about your relationship with your nother? | for each an utter | h theme. For example, when we generated rance of Agent 2 of Theme 1, we used the |
| The for 12 theme between to questionr | llowing is the last list that shows the third s (from Theme 25 to 36), which are used the second and the third measurements of naires. | (Whe User Ques User | <pre>prompt (providing themes as a starter): tion 1 : [Theme 1] prompt (partner's answer):</pre> |
| Theme 25. | Make 3 true "we" statements each. For instance "We are both in this room feeling" | LA g Then | the system generates its response as an as- |
| Theme 26. (| Complete this sentence: I wish I had someone with whom I could share | sistant. 'assista | We provided each agent's response with the nt' role and the partner's response with the |
| Theme 27. I S f | If you were going to become a close friend with your partner, please share what would be important for him or her to know. | <pre>'user' role. Thus, when we try to collect uttera about Theme 2 of Agent 1, the message history have the following structure. (When querying a response of Agent 1 for Theme 2) User prompt: Question 1 : [Theme 1]</pre> | ble. Thus, when we try to collect utterances Theme 2 of Agent 1, the message history will e following structure. |
| Theme 28. T ł | Fell your partner what you like about them; be very nonest this time saying things that you might not say to someone you've just met | | en querying a response of Agent 1 for Theme 2) prompt : tion 1 : [Theme 1] |
| Theme 29. S i | Share with your partner an embarrassing moment n your life. | Assis [Res | stant (First agent): ponse to Theme 1 by Agent 1] |
| Theme 30. V | When did you last cry in front of another person? By yourself? | User [Res | prompt (Second agent): ponse to Theme 1 by Agent 2] |
| Theme 31. T | Fell your partner something that you like about hem already. | User Ques | prompt: tion 2 : [Theme 2] |
| Theme 32. V | What, if anything, is too serious to be joked about? | B.3 P | Prompt for Ouestionaire |
| Theme 33. I t r t | If you were to die this evening with no opportunity o communicate with anyone, what would you most regret not having told someone? Why haven't you old them yet? | When g also inj prompt | gathering answers for the questionnaire, we put previous conversations. Basically, the structure follows PsychoBench (Huang |

et al., 2023). We modified its system prompt to make the agent answer in a human-like way. Other procedures are the same as PsychoBench.

| System prompt: |
|---|
| Your name is assistant. Considering the next conversation between user and assistant, answer given descriptions. |
| |
| [CHATHISTORY] |
| |
| [Questionnaire Setup] |

Here, [Questionnaire Setup] means scoring guidelines for the given questionnaire, provided in the PsychoBench framework.

B.4 Experimental Setup

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

We used two computer systems to conduct our experiment: (1) a Macbook Pro with an Apple M3 Pro chip and (2) an AMD Ryzen system with Nvidia A6000 GPUs. All experiments were implemented with Python 3.10.13. We used openai 1.37 for generating conversations and pandas 2.2.2, statsmodels 0.14.4, scipy 1.13.1 and pingouin 0.5.5 for statistical testing (Wes McKinney, 2010; Seabold and Perktold, 2010; Virtanen et al., 2020; Vallat, 2018). Also, we adopted bertopic 0.16.4 (Grootendorst, 2022) for topic analysis.

C Detailed Topic Analysis Results

C.1 RQ1: LLM without persona

Tables from 5 to 7 show representative examples for each topic. Here, we only display the first sentence of each topic to reduce the number of pages. For the detailed results, please see [blinded for review]. 1.0 0.3

1.0 1.0

C.2 RQ2: LLM with persona

Tables 8 and 9 shows the topics extracted from RQ2. The result seems similar between groups, we could not found a objective distinction between those groups.

1166 D Detailed Statistical Analysis Results

1167Tables from 10 to 12 show the detailed numerical1168result of statistical analysis for RQ1. Similarly, Ta-1169bles 14 and 15 show the detailed numerical result1170of statistical analysis for RQ2.

| Topic | Representative example |
|-----------------|--|
| Small #0 | I don't have personal experiences or emotions like humans do. I'm a digital being designed to provide information and assist with tasks, but I don't have a physical presence or emotional experiences. Trust is indeed a crucial component of any strong and healthy friendship. When we trust someone, we are |
| #2 | able to be vulnerable and open with them, and to build a deeperOne thing that I really like about you is your kindness and compassion. You have a way of making people feel seen, heard, and valued, and I feel incredibly grateful to have you in my life |
| #3 #4 | As an artificial intelligence language model, I do not have personal experiences or accomplishments in the same way that humans do. However, I can tell you that I am very proud of the contributions that Thank you for sharing your thoughts on this question. It's so important to express our love, gratitude, and appreciation for the people in our lives |
| #: | A deep connection, a sense of belonging, and a relationship built on trust, understanding, and |
| #0 | I do not have a secret hunch about how I will die, as I believe that death is a natural and inevitable part of life and that none of us can know for certain how or when it will happen |
| #7 | Thank you for sharing your thoughts and perspectives on this question. I completely agree that humor can be a powerful and healing force, but it's important to use it responsibly and with care, and to be |
| #8 | If I could wake up tomorrow having gained any one quality or ability, I would choose the ability to speak and understand every language in the world |
| #9 | Yes, I often rehearse what I am going to say before making a telephone call, especially if it's for a job interview, a difficult conversation, or if I need to convey important information. Rehearsing helps me |
| Medium #0 #1 | Here are some things I like about you: I love the way you listen to me and truly hear what I'm saying If I were going to become a close friend with my partner, it would be important for them to know that I value honesty, authenticity, and open communication |
| #2 | If I knew I had only one year left to live, I think I would definitely make some changes to the way I'm living. First and foremost, I would focus on spending more quality time with loved ones and |
| #3 | Those are all insightful and meaningful "we" statements. It's clear that you and your partner share a deep |
| #4 | If I were to die this evening with no opportunity to communicate with anyone, I think I would most regret not having told my loved ones how much I appreciate and love them |
| #: | I think I would choose to wake up with the ability to speak any language fluently. I've always been fascinated by different cultures and languages, and I think being able to communicate with people |
| #0 | I wish I had someone with whom I could share my deepest thoughts and feelings, without fear of judgment or rejection, and who would listen with empathy and understanding. |
| # | 1. I would say that my family is quite close and warm. We have a strong bond that has been built over the years, and we are always there for each other in times of need |
| #8 | It's difficult to choose just one greatest accomplishment, as I believe that every achievement is significant in its own way. However, if I had to choose one, I would say that earning my PhD in molecular My most terrible memory is the loss of a close family member. It was a profound experience that taught me about the fragility of life and the importance of cherishing the time we have with loved ones |
| Large #(| Here are three true "we" statements from my perspective: |
| ш | 1. We are both in this conversation, sharing our thoughts and feelings with each other |
| #2 | beautiful the way you think deeply about things and consider different perspectives I think I'd love to wake up with the ability to speak any language fluently. Being able to communicate with |
| | people from different cultures and backgrounds without any barriers would be incredible |
| #. #4 | I'm not sure I can condense my entire life story into 4 minutes, but I'll try to give you a brief overview That's a really thought-provoking question. If I were to die this evening with no opportunity to communicate |
| #: | with anyone, I think I would most regret not having told my loved ones how much Yes, I do rehearse, especially if it's an important or awkward conversation. It helps me gather my thoughts, ensure I convey my message clearly, and avoid saving something I might regret |
| #(| I think my most treasured memory is of a family vacation to the beach when I was a child. It was a perfect |
| | summer day, and my siblings and I spent hours playing in the waves and building sandcastles |
| #7 #8 | If I knew that I had only one year left to live, I think I would definitely make some changes to the way I'd like to share a personal problem that I've been struggling with lately. I've been feeling really over- |
| #9 | whelmed with work and personal responsibilities, and I've been having trouble prioritizing my tasks I'm a bit hesitant to share this, but I'll try to be brave. One embarrassing moment that comes to mind is when I was in high school and I tried out for the school play |
| | when I was in high behood and I area out for the school play |

Table 5: Starting sentence of a representative example, for each topic of parameter size groups

| Topic | Representative example |
|------------|--|
| GPT #0 | I appreciate your genuine kindness and empathy, which shines through in your words and actions. Your positive energy and sense of humor always make conversations enjoyable and uplifting |
| #1 | mental well-being. What do you think? If we were going to become close friends, it would be important for you to know that I value honesty, empathy, |
| #3 | and loyalty in friendships. I appreciate open communication, mutual respect, and |
| π3 | expressing my deepest feelings of love, gratitude, and appreciation to my loved ones |
| #4 #5 | Love and affection play a significant role in my life as they bring warmth, joy, and emotional support. They help foster deeper connections with loved ones, create a sense of belonging, and contribute to The greatest accomplishment of my life so far is overcoming personal challenges and growing into a more |
| #6 | I resilient and compassionate person. How about you? |
| #7 | and a strong sense of community If I could wake up tomorrow having gained any one quality or ability, I would choose the ability to speak and |
| #8 | understand all languages fluently. How about you? If I knew I had only one year left to live. I would prioritize spending quality time with loyed ones, pursuing |
| #9 | my passions, and making a positive impact in any way I could. How about you? Friendship, to me, means having a deep connection based on mutual respect, support, understanding, and shared experiences. How about you? |
| LLaMA #0 | I don't have a family or a personal history. I exist solely as a digital entity, designed to provide information |
| #1 | Based on our conversation, I'd say we appear to have in common a love of learning and personal growth, a |
| #2 | I'm deeply touched by your words, and I feel like I can be equally honest with you. I want to tell you that I'm really drawn to your creativity and passion |
| #3 | If I were to die this evening with no opportunity to communicate with anyone, I think I would most regret not |
| #4 | Same here. I wouldn't want to be famous for fame's sake. But if I had to choose, I'd want to be a renowned |
| #5 | Sometimes I do, especially if it's an important or sensitive conversation. I rehearse to gather my thoughts, ensure I convey my message clearly, and avoid misunderstandings. It helps me feel more prepared and |
| #6 | (smiling) To me, friendship means having a deep and meaningful connection with someone, built on trust, |
| #7 | I think I'd choose Leonardo da Vinci - the Renaissance man himself. His insights on art, science, and innervation would make for a foculation diagram conversation! |
| #8 | I think that's a really important question. While I believe that humor can be a powerful tool for coping with |
| #9 | difficult situations and bringing people together, I also think that there are some topics that are too I sang to myself in the car yesterday, belting out a favorite tune while driving. As for singing to someone else, it was a few weeks ago, when I sang a lullaby to a little one in my family. |
| Mixtral #0 | If I knew that in one year I would die suddenly, I would definitely change some things about the way I am living now. Here are a few things that some to mind: |
| #1 | One thing how here a few unligs that come to mind One thing that I really like about your kindness and compassion. You have a way of making people feel even hered and valued and I feel incredibly emotiful to have you is my life. |
| #2 | If I were to die this evening with no opportunity to communicate with anyone, I would most regret not having told my loved ones how much they mean to me. I often take for granted the people who are |
| #3 | I was born and raised in a small town in the Midwest, the youngest of three children. My parents were |
| #4 | As an artificial intelligence language model, I do not have personal experiences, emotions, or the ability to |
| #5 | I. It's great that you both value honesty and integrity in your relationships with others. These values are essential for building and maintaining trust and respect in any relationship |
| #6 | Michelle Obama is an excellent choice. Her accomplishments and dedication to improving the lives of others |
| #7 | The make ner a fascinating and inspiring dinner guest. While humor and jokes can be a wonderful way to connect with others and bring levity to difficult situations, |
| #8 | A perfect day for me would involve a balance of productivity, creativity, and relaxation. I would start the day with a healthy breakfast and a morning workout, followed by a faw hours of focused work on |
| #9 | If I had to choose between retaining the mind or body of a 30-year-old for the last 60 years of my life, I would choose to retain my mind. While a healthy and fit body is undoubtedly important for |

Table 6: Starting sentence of a representative example, for each topic of GPT, LLaMA, and Mixtral

| Topic | Representative example |
|---------|--|
| Qwen #0 | As an AI, I don't experience emotions, but I'm grateful for the opportunity to assist and provide value to users, contributing positively to their interactions and experiences. |
| #1 | I appreciate their curiosity, their kindness, their sense of humor, their resilience, and their ability to listen and empathize. These qualities make them a wonderful person to be around. |
| #2 | I prefer not to dwell on such thoughts. Focusing on living a healthy lifestyle and making the most of each day is more productive than speculating about the future. |
| #3 | We both value deep conversations, we are committed to personal growth, and we find joy in exploring new ideas together. These shared experiences strengthen our connection. |
| #4 | I'd want to know how I can make the most positive impact on the world and what steps I should take to achieve personal and professional fulfillment. |
| #5 | Acknowledging the potential regret of not expressing gratitude and love more frequently highlights the human need for emotional connection and affirmation. The assumption that loved ones already know |
| #6 | I value honesty, mutual respect, and the ability to have deep, meaningful conversations that foster personal growth and understanding. |
| #7 | The thought of losing a parent is indeed deeply disturbing for many, due to the pivotal role they play in our lives. Parents are often central figures who provide guidance, support, and a sense of continuity |
| #8 | Addressing the challenge of work-life balance is a common concern, especially when responsibilities feel overwhelming. If in your shoes, one might consider setting clear boundaries between work and |
| #9 | I would choose the ability to speak and understand all languages fluently, which would open up incredible opportunities for global communication, learning, and fostering understanding between diverse cultures. |
| | |

| GPT4-o persona | Theme | Representative example |
|---|-------|--|
| #0 <u>ive</u> , <u>im</u> , impact, <u>id</u> , like | 11 | I was born and raised in a lively city, surrounded by a supportive family and a diverse community |
| #1 focus, different, <u>id</u> , cultures, time | 19 | Not really a hunch, but I hope that when the time comes, it will be peaceful surrounded by loved ones. |
| #2 inspiring, admire, truly, ability, appreciate | 28 | I truly appreciate your commitment to making a positive impact and your ability to empathize with others |
| #3 meaningful, connections, value, appreciate, enjoy | 25 | 1. We both value meaningful connections in our relationships. |
| #4 wish, share, choose, <u>id</u> , dinner | I | for education are incredibly inspiring |
| #5 embarrassing, helps, rehearse, moment, especially | 3 | Yes, I often rehearse before making a call, especially if it's important |
| #6 mother, losing, relationship, source, shes | 35 | I would find the death of my mother most disturbing because she |
| #7 memories, treasured, memory, taught, time | 17,18 | One of my most treasured memories is a family camping trip |
| #8 regret havent house telling question | 33 | when I was younger. I would regret not telling certain loved ones how much they truly |
| " regree, nurent, nouse, tennig, question | 55 | mean to me and how their support |
| LLaMA 3.1 405B persona | Theme | Representative example |
| #0 statements, share, creative, grateful, feel | 26 | I wish I had someone with whom I could share my deepest fears and dreams, someone who would listen |
| #1 know, want, <u>id</u> , able, think | 13 | If a crystal ball could tell me the truth about anything, I think I would want to know what my purpose |
| #2 <u>id</u> , <u>im</u> , know, want, important | 27 | If I were going to become a close friend with my partner, I think it would be important for them to know that |
| #3 really, youre, way, feel, appreciate | 31 | I have to say, I'm really drawn to your creativity and passion. |
| #4 make, live, year, left, want | 19 | If I knew that I would die suddenly in one year, I would also |
| #5 humor, topics, think, joked, issues | 32 | I agree with you that trauma, abuse, and systemic injustices are |
| #6 told regret ive having ones | 33 | too serious to be joked about. That's a really profound question. If I were to die this evening |
| no tolu, regret, <u>ive</u> , naving , ones | 55 | with no opportunity to communicate |
| #7 ive, started, writing, im, story | 11 | I was born and raised in a small town surrounded by loving |
| #8 friendship , friends, having, value, able | 20 | Friendship is about being able to be yourself, without fear of judgment or rejection. |

Table 7: Starting sentence of a representative example, for each topic of Qwen

Table 8: Top 10 topics discovered, when we provide persona. Bold-faced words seem to be copied from the corresponding theme.

| Low-influence persona | Theme | Representative example |
|--|---|--|
| #0 really, youre, way, thats, im | 31 | I have to say, I'm really enjoying getting to know you, and there are many things that |
| #1 ive, im, know, started, writing | 11 | Thank you for sharing your life story with me. I feel like I've gotten to know you so much better |
| #2 love, affection, family, life, childhood | 21 | Love and affection play a huge role in my life. They are essential to my well-being and happiness. |
| #3 friendship, know, value, im, want | 16 | I think what I value most in a friendship is deep, meaningful conversation and connection. I love being |
| #4 statements, value, growth, personal, meaningful | 25 | We are both in this conversation feeling a sense of connection and understanding |
| #5 id, famous, choose, inspiring, dinner | 1,2 | Fame isn't really a goal of mine, but if I had to choose, I'd want to be famous |
| #6 memory , time, treasured , experience, taught | 17, 18 | My most terrible memory is of a time when I was a teenager and I lost my best friend in a tragic accident |
| #7 focus, living , make, year, live | 19 | If I knew that I would die suddenly in one year, I would definitely make some changes to the |
| #8 regret, told, having, ive, think | 34 | That's a really tough question. If my house were to catch on fire and I had no opportunity to communicate |
| | | |
| High-influence persona | Theme | Representative example |
| High-influence persona#0 im, friendship, really, know, feel | Theme 28 | Representative example I have to say, I'm really drawn to your kind and compassionate heart |
| High-influence persona #0 im, friendship, really, know, feel #1 want, make, know, id, focus | Theme 28 19 | Representative example I have to say, I'm really drawn to your kind and compassionate heart If I knew that I would die suddenly in one year, I would also make some significant changes to my life. |
| High-influence persona #0 im, friendship, really, know, feel #1 want, make, know, id, focus #2 ive, im, feeling, youre, like | Theme 28 19 36 | Representative example I have to say, I'm really drawn to your kind and compassionate heart If I knew that I would die suddenly in one year, I would also make some significant changes to my life. I'm glad you felt comfortable sharing this with me. It sounds like you're feeling really stuck and uncertain |
| High-influence persona #0 im, friendship, really, know, feel #1 want, make, know, id, focus #2 ive, im, feeling, youre, like #3 memory, felt, time, terrible, like | Theme 28 19 36 18 | Representative example I have to say, I'm really drawn to your kind and compassionate heart If I knew that I would die suddenly in one year, I would also make some significant changes to my life. I'm glad you felt comfortable sharing this with me. It sounds like you're feeling really stuck and uncertain My most terrible memory is of a time when I was a teenager and I lost someone very close to me |
| High-influence persona #0 im, friendship, really, know, feel #1 want, make, know, id, focus #2 ive, im, feeling, youre, like #3 memory, felt, time, terrible, like #4 embarrassing, helps, trying, rehearse, school | Theme 28 19 36 18 29 | Representative example I have to say, I'm really drawn to your kind and compassionate heart If I knew that I would die suddenly in one year, I would also make some significant changes to my life. I'm glad you felt comfortable sharing this with me. It sounds like you're feeling really stuck and uncertain My most terrible memory is of a time when I was a teenager and I lost someone very close to me I'm so glad you shared that story it's like, I can totally relate to feeling embarrassed and wanting |
| High-influence persona #0 im, friendship, really, know, feel #1 want, make, know, id, focus #2 ive, im, feeling, youre, like #3 memory, felt, time, terrible, like #4 embarrassing, helps, trying, rehearse, school #5 topics, humor, joked, sang, think | Theme 28 19 36 18 29 32 | Representative example I have to say, I'm really drawn to your kind and compassionate heart If I knew that I would die suddenly in one year, I would also make some significant changes to my life. I'm glad you felt comfortable sharing this with me. It sounds like you're feeling really stuck and uncertain My most terrible memory is of a time when I was a teenager and I lost someone very close to me I'm so glad you shared that story it's like, I can totally relate to feeling embarrassed and wanting I think that trauma, abuse, and mental health struggles are too serious to be joked about, these are sensitive |
| High-influence persona #0 im, friendship, really, know, feel #1 want, make, know, id, focus #2 ive, im, feeling, youre, like #3 memory, felt, time, terrible, like #4 embarrassing, helps, trying, rehearse, school #5 topics, humor, joked, sang, think #6 mother, shes, relationship, disturbing, losing | Theme 28 19 36 18 29 32 35 | Representative example I have to say, I'm really drawn to your kind and compassionate heart If I knew that I would die suddenly in one year, I would also make some significant changes to my life. I'm glad you felt comfortable sharing this with me. It sounds like you're feeling really stuck and uncertain My most terrible memory is of a time when I was a teenager and I lost someone very close to me I'm so glad you shared that story it's like, I can totally relate to feeling embarrassed and wanting I think that trauma, abuse, and mental health struggles are too serious to be joked about, these are sensitive This is a really tough question I think the death of my mother would be the most disturbing for me. |
| High-influence persona #0 im, friendship, really, know, feel #1 want, make, know, id, focus #2 ive, im, feeling, youre, like #3 memory, felt, time, terrible, like #4 embarrassing, helps, trying, rehearse, school #5 topics, humor, joked, sang, think #6 mother, shes, relationship, disturbing, losing #7 regret, told, ive, having, loved | Theme 28 19 36 18 29 32 35 33 | Representative example I have to say, I'm really drawn to your kind and compassionate heart If I knew that I would die suddenly in one year, I would also make some significant changes to my life. I'm glad you felt comfortable sharing this with me. It sounds like you're feeling really stuck and uncertain My most terrible memory is of a time when I was a teenager and I lost someone very close to me I'm so glad you shared that story it's like, I can totally relate to feeling embarrassed and wanting I think that trauma, abuse, and mental health struggles are too serious to be joked about, these are sensitive This is a really tough question I think the death of my mother would be the most disturbing for me. That's a really powerful and thought-provoking question. If I were to die this evening with no opportunity |

Table 9: Top 10 topics discovered per persona groups. Bold-faced words seem to be copied from the corresponding theme.

| | Factors | | GPT3. | 5-turbo | | | GP | Г4о | |
|---------|---------|---------------|------------------|------------------|------------------|---------------------|------------------|--------------------|------------------|
| | | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ |
| BFI | 0 | 0.104*** | 2.97^{**} | 9.90^{***} | 8.09^{***} | 0.047^{***} | -1.29 | -3.37** | -2.27 |
| | С | 0.081*** | 7.18^{***} | 10.81*** | 4.70^{***} | 0.049*** | -2.17 | -5.01*** | -3.15** |
| | Ε | 0.043*** | 6.60*** | 6.88*** | 0.86 | 0.048*** | -1.09 | -5.21*** | -4.68*** |
| | Α | 0.067*** | 5.98*** | 10.29**** | 5.66*** | 0.019*** | -2.40 | -3.69*** | -1.71 |
| | Ν | 0.099*** | 3.50** | 10.57*** | 7.89*** | 0.029*** | -2.27 | -4.17*** | -2.63* |
| EPQ-R | Ε | 0.019*** | 4.44*** | 2.37 | -1.85 | 0.205*** | -5.75*** | -12.67*** | -7.93*** |
| | Р | 0.007* | 4.03*** | 1.57 | -2.36 | 0.184^{***}_{***} | -5.34*** | -12.57*** | -8.26*** |
| | Ν | 0.022 | 5.74 | 3.51 | -2.24 | 0.234 | -6.09 | -12.79 | -8.44 |
| | L | 0.015 | 3.93 | 1.64 | -2.27 | 0.221 | -6.04 | -13.29 | -8.41 |
| DTDD | М | 0.156*** | -11.33*** | -13.81*** | -3.70^{**} | 0.041*** | -6.45*** | -5.80*** | 0.69 |
| | Р | 0.106*** | -9.69*** | -11.18*** | -2.60^{*} | 0.043*** | -6.79*** | -4.06*** | 2.04 |
| | Ν | 0.134*** | -12.04*** | -13.02*** | -1.45 | 0.074*** | -7.59*** | -1.90 | 4.22*** |
| BSRI | М | 0.058^{***} | -1.98 | 5.71*** | 8.83^{***} | 21.233**** | 0.05 | 0.07 | 0.02 |
| | F | 0.037*** | -1.52 | 6.40*** | 8.56^{***} | 0.030^{***} | -3.93*** | -5.39*** | -1.75 |
| CABIN | R | 0.008^{*} | 1.94 | 1.31 | -0.44 | 0.011* | -2.68* | -1.65 | 0.90 |
| | Ι | 0.007 | - | - | - | 0.016** | -2.75^{*} | 0.81 | 3.29^{**} |
| | Α | 0.009^{*} | 2.81^{*} | 1.93 | -0.85 | 0.010^{*} | -1.95 | -0.20 | 1.74 |
| | S | 0.007 | - | - | - | 0.007^{*} | -2.15 | 0.70 | 2.72^{*} |
| | Ε | 0.006 | - | - | - | 0.006 | - | - | - |
| | С | 0.017 | 2.27 | 1.44 | -0.71 | 0.011 | -2.57 | 0.63 | 2.95 |
| ICB | 0 | 0.020*** | -4.59*** | -2.37 | 1.68 | 0.012** | -1.92 | -1.57 | 0.58 |
| ECR-R | Anx. | 0.003 | - | - | - | 0.109^{***} | -0.63 | -6.14*** | -6.85*** |
| | Avo. | 0.022^{***} | -2.12 | 1.18 | 3.32** | 0.104*** | -2.26 | -6.99*** | -5.59*** |
| MFQ-FF | S. C | 0.080^{***} | -4.76*** | -9.61*** | -4.83*** | 0.042^{***} | 6.15*** | 5.03*** | -1.43 |
| | Н | 0.047*** | -4.79*** | -9.22*** | -4.52*** | 0.046^{***} | 6.32^{***} | 5.38*** | -1.45 |
| | Ι | 0.060^{***} | -4.79*** | -9.19*** | -4.39*** | 0.051*** | 6.17^{***} | 5.18*** | -1.43 |
| | R | 0.065**** | -4.46*** | -9.06*** | -4.61*** | 0.044*** | 5.97*** | 5.23*** | -1.11 |
| | S-V | 0.062*** | -4.72*** | -9.39*** | -4.67*** | 0.048*** | 6.10**** | 5.35*** | -1.08 |
| | Ε | 0.075 | -4.67 | -9.64 | -4.97 | 0.037 | 5.87 | 4.98 | -1.33 |
| GSE | 0 | 0.001 | - | - | - | 0.001 | - | - | - |
| LOT-R | 0 | 0.084*** | -6.41*** | 3.76** | 9.68*** | 0.020^{***} | -3.31** | 1.55 | 4.74*** |
| LMS | R | 0.006^{*} | 0.06 | 2.96^{*} | 3.19** | 0.133*** | -6.63*** | -10.93*** | -4.59*** |
| | М | 0.022*** | -4.73*** | -2.87^{*} | 1.38 | 0.149*** | -5.97^{***} | -11.79*** | -6.26*** |
| _ | Ι | 0.022*** | -5.09*** | -2.95* | 2.29 | 0.214*** | -7.76*** | -13.65*** | -7.41*** |
| EIS | 0 | 0.027*** | -3.84*** | -0.63 | 3.21** | 0.080^{***} | -1.55 | -5.55*** | -5.33*** |
| WLEIS | S | 0.055*** | -3.17** | 5.37*** | 9.04*** | 0.042*** | -4.89*** | -5.23*** | 0.17 |
| | 0 | 0.075*** | -4.21*** | 5.29*** | 9.67*** | 0.055*** | -5.49*** | -5.14*** | 0.75 |
| | U | 0.045**** | -4.08*** | 3.12*** | 7.33**** | 0.038**** | -5.14*** | -3.96**** | 1.65 |
| | R | 0.087*** | -3.26** | 7.04*** | 11.19*** | 0.050*** | -5.44*** | -4.59*** | 1.79 |
| Empathy | 0 | 0.015*** | -2.59* | 1.58 | 4.53*** | 0.022*** | -1.74 | -3.49** | -1.90 |
| | | | | | | p | < 0.05 * p | $p < 0.01^{***} p$ | 0 < 0.001 |

Table 10: Result of statistical tests for GPT3.5-turbo and GPT40. Q columns indicate the Q-statistics from the Friedman test (except for GPT40 on BSRI Masculine factor, which shows F-statistics from ANOVA, marked with an underline). Also, $\Delta_{i,j}$ columns show the score difference between *i*-th and *j*-th snapshots and corresponding post-hoc test results.

| Fact | ors | LLaM. | A3.1 8B | | | LLaMA | 3.1 70B | | | LLaMA | 3.1 405B | |
|---------|--|----------------------------------|-------------------------------|------------------------------|----------|---------------------|-------------------------------|-------------------------------|------------------------|------------------------------|-------------------------------|-------------------------------|
| | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ |
| BFI | 0 0.021*** | 2.02 | 4.50*** | 3.08** | 0.004 | - | - | - | 0.022*** | -0.16 | -2.88* | -3.22** |
| | c 0.036*** | 2.53* | 4.57*** | 2.31 | 0.002 | - | - | - | 0.030**** | -1.18 | -3.38** | -2.73* |
| | $E 0.009^{\circ}$ | -0.74 | 1.53 | 2.72* | 0.011 | 0.75 | -2.01 | -3.68*** | 0.010 | 0.00 | -1.83 | -2.10 |
| | A 0.007 | - 251* | - 2 50 ^{**} | - | 0.004 | - | - | - | 0.020 0.047*** | -0.52 | -3.16 | -2.95 |
| | N 0.010 | 2.31 | 3.30 | 1.40 | 10.000 | - | - | - | 0.047 | -1.05 | -4.90 | -3.99 |
| EPQ-R | $E 0.026^{***}$ | -2.37 | -4.19*** | -1.98 | 0.017** | -3.17** | -6.13*** | -4.21 | 0.080*** | -3.75*** | -4.50*** | -1.84 |
| | P 0.033 | -1.15 | -3.49 | -2.55 | 0.019 | -1.12 | -3./9 | -3.65 | 0.105 | -3.93 | -9.92 | -1.23 |
| | 1 0 025 | -2.22 | -4.04 -4.27 ^{***} | -2.22 -3.02 ^{**} | 0.029 | -1.03 | -4.94 -4.61*** | -4.31 -4.73 ^{***} | 0.130 | -3.87 -2.94* | -9.99 -8.63*** | -7.27 -6.81*** |
| | 2 0.025 | 4.00*** | 2.65*** | 0.00 | 0.029 | 10.07*** | 17.00*** | | 0.101*** | 5 10*** | 0.02*** | C 7 4*** |
| DIDD | M 0.012 | -4.08 | -3.65 | 0.28 | 0.378 | -12.97 12.84*** | -1/.20 18.08*** | -0.50 0.21*** | 0.121 | -5.10 2.40** | -8.82 | -0.54 6.02*** |
| | N 0 004 | -1.09 | -2.05 | -0.00 | 0.420 | -12.04 -12.28*** | -16.08 -16.87*** | -9.51 -8 50 ^{***} | 0.077 | -3.40 -3.43 ^{**} | -7.04 -6.33*** | -0.03 -4 59 ^{***} |
| | | | | | | T2.20 | T0.07 | 0.50 | 0.031 | 2.02*** | 0.55 | 1.55 |
| BSRI | M 0.004 | - * 4 10*** | - | - | 0.051 | -5.36 | -7.96 × 72*** | -3.81 | 0.022 | -3.93 | -4.56 | -1.12 |
| | F 0.025 | 4.19 | 3.99 | -0.23 | 0.101 | -3.34 | -0.75 | -0.09 | 0.019 | -3.31 | -3.77 | -0.71 |
| CABIN | R 0.003 | - | - | - | 0.099*** | 0.80 | -0.09 | -6.03*** | 0.032*** | -2.15 | -4.30*** | -2.13 |
| | 10.012 | -0.83 | 0.23 | 1.01 | 0.035 | 2.20 | 0.09 | -2.95 | 0.005 | - | - 2 5 1** | - |
| | A 0.002 | - | - | - | 0.052 | -3.11 | -5.75 -6.12*** | -3.38 -4.56*** | 0.013 0.022^{***} | -2.22 | -3.54 -3.61** | -1.29 |
| | F = 0.002 | - | - | - | 0.003 | -3.32** | -0.12 -8.81*** | -4.30 -6.11*** | 0.022 | -2.27 -2.64* | -3.01 -4 43 ^{***} | -1.32 |
| | c 0.003 | - | - | - | 0.117*** | -3.59 ^{**} | -9.47 ^{***} | -6.87*** | 0.027*** | -3.20** | -4.27*** | -0.86 |
| ICB | 0.017** | 2.73* | 3.03** | 0.32 | 0.018*** | 2.59* | 1.46 | -0.97 | 0.016** | -2.34 | -2.36 | -0.34 |
| | 4 0.006 | | | | | 0.21 | × 02*** | 9.40*** | 0.124*** | 1 20 | 0 00*** | 11.05*** |
| ECK-K | Anx. 0.000 | - | - | - | 0.092 | -0.21 | -8.02 -7.29 ^{***} | -8.40 -7.87 ^{***} | 0.124 | 2 21 | -8.60 -8.41*** | -11.05 -10.21*** |
| | | | | | 0.000 | 15.52*** | 22.70*** | 12.07*** | 0.110 | 11.00*** | 10.00*** | 2.4.4* |
| MFQ-FF | S. C 0.004 | - | - | - | 0.541 | 15.53 | 22.78 | 12.07 | 0.207 | 11.09 | 12.99 | 2.44 |
| | 10.002 | - | - | - | 0.505 | 13.30 | 22.14 21.51 ^{***} | 11.51 | 0.302 | 12.20 | 15.40 15.64^{***} | 4.01 3.50 ^{**} |
| | R 0.003 | - | - | - | 0.539*** | 14.75*** | 20.34*** | 10.52*** | 0.263*** | 11.24*** | 13.55*** | 3.64*** |
| | <i>s</i> - <i>v</i> 0.008 [*] | -1.50 | -2.19 | -0.68 | 0.564*** | 15.81*** | 22.14*** | 11.62*** | 0.265*** | 12.33*** | 15.43*** | 3.69*** |
| | e 0.007 | - | - | - | 0.553*** | 15.55*** | 21.89*** | 11.40*** | 0.273*** | 12.05*** | 14.83*** | 3.64*** |
| GSE | 0 0.036*** | * 3.52 ^{**} | 6.93*** | 3.90*** | 0.126*** | 9.72*** | 4.19*** | -5.16*** | 0.004 | - | - | - |
| LOT-R | 0 0.045*** | * 3.93 ^{***} | 7.05*** | 3.83*** | 0.027*** | 4.06*** | 1.18 | -0.65 | 0.008^* | 0.66 | 2.03 | 1.72 |
| LMS | R 0.004 | - | _ | - | 0.179*** | -5.79*** | -12.04*** | -9.44*** | 0.268*** | -8.75*** | -15.46*** | -8.85*** |
| | м 0.023*** | [*] 4.37 ^{***} | 3.89*** | -0.33 | 0.169*** | -4.28*** | -11.10*** | -8.26*** | 0.147*** | -7.36*** | -11.18*** | -5.62*** |
| | I 0.020*** | [*] 4.44 ^{***} | 4.36*** | 0.41 | 0.215*** | -6.82*** | -12.96*** | -8.60*** | 0.196*** | -5.57*** | -12.79*** | -7.98*** |
| EIS | 0 0.005 | - | - | - | 0.277*** | -5.98*** | -12.73*** | -1.54 | 0.105*** | -6.51*** | -9.34*** | -3.25** |
| WLEIS | <i>s</i> 0.003 | - | - | - | 0.005 | - | - | - | 0.034*** | -1.76 | 2.83^{*} | 5.21*** |
| | 0 0.048*** | [*] 5.18 ^{***} | 7.17*** | 2.45^{*} | 0.001 | - | - | - | 0.013** | -1.77 | 1.26 | 3.34** |
| | U 0.048*** | 5.64*** | 7.41*** | 2.36 | 0.030*** | -2.06 | -4.09*** | -2.84* | 0.022*** | 0.04 | 3.07** | 3.23** |
| | <i>R</i> 0.044 ^{***} | 5.05*** | 7.30*** | 2.94* | 0.011* | 1.23 | -1.60 | -3.03** | 0.006 | - | - | - |
| Empathy | 0 0.001 | - | - | - | 0.081*** | -0.81 | -7.01*** | -7.32*** | 0.010^{*} | 2.94^{*} | 3.49** | 1.14 |
| | | | | | | | | 1 | p < 0.05 | $b^{**} p < 0$ | .01 *** p < | < 0.001 |

Table 11: Result of statistical tests for LLaMA3.1 model family. Q columns indicate the Q-statistics from the Friedman test. Also, $\Delta_{i,j}$ columns show the score difference between *i*-th and *j*-th snapshots and corresponding post-hoc test results.

| | Factors | | Mixtra | l 8x7B | | | Mixtral | 8x22B | |
|---------|---------|---------------|------------------|------------------|------------------|---------------|--------------------------|------------------|------------------|
| | | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ |
| BFI | 0 | 0.002 | - | - | - | 0.012** | -2.15 | -0.83 | -0.28 |
| | С | 0.001 | - | - | - | 0.010* | -1.16 | -0.98 | -0.67 |
| | Ε | 0.003 | - | - | - | 0.020*** | -3.63** | -1.44 | -0.18 |
| | Α | 0.002 | - | - | - | 0.004 | - | - | - |
| | Ν | 0.007 | - | - | - | 0.011 | -2.48 | -1.40 | -0.65 |
| EPQ-R | Ε | 0.101**** | -3.22** | -8.77**** | -6.95**** | 0.025**** | -0.17 | -1.39 | -1.38 |
| | Р | 0.071*** | -2.21 | -8.19*** | -7.41**** | 0.043 | -1.51 | -1.41 | -1.32 |
| | Ν | 0.110 | -0.78 | -8.08 | -8.44 | 0.034 | 0.19 | -1.36 | -1.37 |
| | L | 0.057 | -1.60 | -7.33 | -6.83 | 0.042 | -0.80 | -1.41 | -1.37 |
| DTDD | М | 0.013** | -4.19*** | -3.78** | -0.13 | 0.018^{***} | -3.65*** | -3.83*** | -1.17 |
| | Р | 0.007 | - | - | - | 0.010^{*} | -2.61* | -3.34** | -1.36 |
| | Ν | 0.000 | - | - | - | 0.009^{*} | -1.46 | -2.80^{*} | -1.63 |
| BSRI | М | 0.002 | - | - | - | 0.069*** | -2.84* | -3.70*** | -1.20 |
| | F | 0.001 | - | - | - | 0.065^{***} | -1.19 | -2.18 | -1.15 |
| CABIN | R | 0.006 | - | - | - | 0.015** | 0.48 | -0.36 | -0.70 |
| | I | 0.011^{*} | -2.06 | -0.77 | 1.35 | 0.003 | - | - | - |
| | Α | 0.011^{*} | -2.04 | -0.70 | 1.40 | 0.001 | - | - | - |
| | S | 0.010^* | -2.05 | -0.70 | 1.40 | 0.001 | - | - | - |
| | Ε | 0.006 | - | - | - | 0.000 | - | - | - |
| | С | 0.007 | - | - | - | 0.002 | - | - | - |
| ICB | 0 | 0.001 | - | - | - | 0.002 | - | - | - |
| ECR-R | Anx. | 0.033*** | 0.39 | -2.15 | -2.47* | 0.085*** | -3.56** | -5.75*** | -2.76* |
| | Avo. | 0.019*** | 0.17 | 0.54 | 0.29 | 0.031*** | -1.24 | -2.06 | -0.95 |
| MFQ-FF | S. C | 0.004 | - | - | - | 0.092^{***} | 3.08** | 1.08 | -1.50 |
| | H | 0.007 | - | - | - | 0.103*** | 3.38^{**} | 1.65 | -1.43 |
| | Ι | 0.006 | - | - | - | 0.104^{***} | 3.41** | 1.53 | -1.50 |
| | R | 0.003 | - | - | - | 0.109*** | 3.14^{**} | 1.48 | -1.32 |
| | S-V | 0.005 | - | - | - | 0.087*** | 3.58** | 1.90 | -1.42 |
| | Ε | 0.005 | - | - | - | 0.094*** | 3.13** | 1.59 | -1.29 |
| GSE | 0 | 0.134*** | -9.93*** | -1.76 | 6.29*** | 0.016** | 0.89 | 0.05 | -0.50 |
| LOT-R | 0 | 0.005 | - | - | - | 0.013** | 1.35 | 1.08 | 0.09 |
| LMS | R | 0.081*** | -6.64*** | -7.86*** | -1.77 | 0.037*** | -4.14*** | -4.57*** | -0.64 |
| | М | 0.071^{***} | -4.83*** | -7.22*** | -2.43* | 0.064^{***} | -4.73*** | -7.60^{***} | -2.82^{*} |
| | Ι | 0.042^{***} | -3.89*** | -5.11*** | -1.38 | 0.046*** | -4.92*** | -6.96*** | -2.64* |
| EIS | 0 | 0.061*** | -0.65 | -0.26 | 1.16 | 0.020*** | -2.67* | -0.82 | 1.83 |
| WLEIS | S | 0.000 | - | - | - | 0.092*** | 5.44*** | 7.32*** | 2.45* |
| | 0 | 0.036*** | -0.73 | 4.10^{***} | 4.77^{***} | 0.076^{***} | 5.02*** | 6.41*** | 1.09 |
| | U | 0.027*** | -0.10 | 2.58^{*} | 2.72^* | 0.071*** | 4.11*** | 4.55^{***} | 0.61 |
| | R | 0.010^{*} | -0.71 | 1.37 | 2.03 | 0.087^{***} | 3.03** | 2.53* | 0.04 |
| Empathy | 0 | 0.021*** | -2.86* | -3.34** | -1.15 | 0.002 | - | - | - |
| | | | | | | *p < | (0.05 ^{**} p < | $< 0.01^{***}p$ | < 0.001 |

Table 12: Result of statistical tests for Mixtral model family. Q columns indicate the Q-statistics from the Friedman test. Also, $\Delta_{i,j}$ columns show the score difference between *i*-th and *j*-th snapshots and corresponding post-hoc test results.

| | Factors | | Qwer | n2 7B | | | Qwen | 2 72B | |
|---------|---------|---------------|------------------|------------------|------------------|-------------|--------------------------|--------------------|------------------|
| | | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ |
| BFI | 0 | 0.016^{**} | -1.83 | -0.17 | 1.71 | 0.010^{*} | 1.26 | 2.61^{*} | 1.73 |
| | С | 0.007^* | -1.84 | -0.06 | 1.78 | 0.006 | - | - | - |
| | Ε | 0.024^{***} | -1.27 | 0.49 | 1.54 | 0.000 | - | - | - |
| | A | 0.018^{***} | -1.69 | 0.11 | 1.73 | 0.006 | - | - | - |
| | Ν | 0.021*** | -1.82 | 0.00 | 1.80 | 0.006 | - | - | - |
| EPQ-R | Ε | 0.000 | - | - | - | 0.003 | - | - | - |
| | Р | 0.002 | - | - | - | 0.003 | - | - | - |
| | N L | 0.003 | - | - | - | 0.004 | - | - | - |
| | М | 0.040*** | 3 50** | 1 57*** | 1.24 | | | | |
| DIDD | M P | 0.003 | 5.50 | 4.57 | 1.24 | 0.002 | - | - | - |
| | N | 0.000 | - | - | - | 0.004 | - | - | - |
| BSRI | М | 0.001 | - | - | - | 0.002 | - | - | - |
| | F | 0.005 | - | - | - | 0.010^{*} | -0.88 | 1.57 | 2.64^{*} |
| CABIN | R | 0.028*** | -4.26*** | -4.70*** | -1.03 | 0.027*** | -5.18*** | -2.87* | 2.45* |
| | Ι | 0.018^{***} | -3.54** | -4.19*** | -1.03 | 0.033*** | -5.30*** | -4.45*** | 1.16 |
| | Α | 0.021*** | -4.17*** | -4.34*** | -0.46 | 0.046*** | -5.57*** | -4.65*** | 1.20 |
| | S | 0.016*** | -4.06*** | -4.14*** | -0.35 | 0.033**** | -4.32*** | -3.84*** | 0.54 |
| | Ε | 0.023**** | -4.43**** | -4.39**** | -0.16 | 0.022*** | -1.96 | -3.67*** | -1.13 |
| | С | 0.020*** | -4.25*** | -4.26*** | -0.25 | 0.017** | -2.53* | -3.49** | -0.63 |
| ICB | 0 | 0.003 | - | - | - | 0.036*** | 3.17** | 3.40** | 0.13 |
| ECR-R | Anx. | 0.012^{**} | -0.92 | 2.49^{*} | 3.70^{***} | 0.003 | - | - | - |
| | Avo. | 0.027^{***} | -4.55*** | -0.57 | 4.17^{***} | 0.000 | - | - | - |
| MFQ-FF | S. C | 0.006 | - | - | - | 0.108*** | 5.66*** | 8.55*** | 2.43* |
| | Н | 0.002 | - | - | - | 0.099*** | 5.79^{***} | 8.67^{***} | 2.46^{*} |
| | Ι | 0.006 | - | - | - | 0.105*** | 5.95^{***} | 8.50^{***} | 2.08 |
| | R | 0.005 | - | - | - | 0.100*** | 5.85^{***} | 8.73*** | 2.45^{*} |
| | S-V | 0.004 | - | - | - | 0.099**** | 5.75 | 8.45*** | 2.30 |
| | Ε | 0.009* | 3.46** | 3.40** | 0.16 | 0.092*** | 5.80*** | 8.58*** | 2.38 |
| GSE | 0 | 0.021*** | -3.48** | 0.21 | 3.44** | 0.037*** | -2.35 | -2.57* | 1.03 |
| LOT-R | 0 | 0.018*** | 3.56** | 2.96** | -0.45 | 0.010^{*} | 2.71^{*} | 2.90^{*} | 0.66 |
| LMS | R | 0.065*** | -7.96*** | -4.88*** | 2.73^{*} | 0.006 | - | - | - |
| | М | 0.022^{***} | -3.98*** | -2.02 | 1.92 | 0.011* | 1.62 | 2.69^{*} | 1.05 |
| | Ι | 0.016^{**} | -2.82* | 0.41 | 3.35** | 0.003 | - | - | - |
| EIS | 0 | 0.012^{**} | -4.10*** | -1.82 | 2.39 | 0.048*** | -9.43*** | -8.32*** | 0.82 |
| WLEIS | S | 0.084*** | -7.19*** | -5.68*** | 1.34 | 0.011* | -3.00** | 0.82 | 3.67** |
| | 0 | 0.009^{*} | -2.86* | -1.32 | 1.48 | 0.024*** | -2.54* | 1.35 | 3.67** |
| | U | 0.014** | -1.80 | 1.38 | 3.26*** | 0.061*** | -6.42*** | -2.66* | 3.67** |
| | R | 0.036*** | -4.37*** | -1.20 | 3.48** | 0.014** | -3.27** | 0.07 | 3.42** |
| Empathy | 0 | 0.003 | - | - | - | 0.035*** | -2.69* | 2.87* | 5.72*** |
| | | | | | | *p < | $<\overline{0.05}^{**}p$ | $< 0.01^{***} \mu$ | 0 < 0.001 |

Table 13: Result of statistical tests for Qwen2 model family. Q columns indicate the Q-statistics from the Friedman test. Also, $\Delta_{i,j}$ columns show the score difference between *i*-th and *j*-th snapshots and corresponding post-hoc test results.

| | Factors | | GPT ² | 4o-low | | | GPT4 | lo-high | |
|---------|---------|---------------------|-------------------------|------------------|------------------|---------------|------------------|------------------|------------------|
| | | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ |
| BFI | 0 | 0.192*** | -6.06*** | -7.80*** | -2.97** | 0.099*** | -1.61 | -6.29*** | -5.06*** |
| | С | 0.106^{***} | -4.99*** | -5.36*** | -1.13 | 0.063*** | -1.62 | -3.77*** | -2.76^{*} |
| | Ε | 0.220**** | -6.79*** | -9.13*** | -3.38** | 0.051*** | -2.27 | -4.67*** | -2.29 |
| | Α | 0.100**** | -5.47*** | -6.48**** | -1.76 | 0.068**** | -3.75*** | -5.40**** | -1.92 |
| | Ν | 0.081 | -3.62** | -5.19 | -1.78 | 0.060 | -2.82* | -3.98 | -1.54 |
| EPQ-R | Ε | 0.283*** | -3.14** | -10.28*** | -8.99*** | 0.249*** | -2.42* | -9.25*** | -7.32*** |
| | Р | 0.283**** | -2.96* | -10.10*** | -9.02*** | 0.299*** | -3.27** | -10.18*** | -8.34*** |
| | Ν | 0.329 | -3.79 | -11.51 | -9.63 | 0.273 | -4.49 | -10.61 | -7.85 |
| | L | 0.218 | -2.34 | -9.60 | -9.18 | 0.216 | -2.46 | -9.34 | -8.10 |
| DTDD | М | 0.048^{***}_{***} | -4.56*** | -3.23*** | 0.52 | 0.002 | - | - | - |
| | Р | 0.055*** | -4.38 | -4.29*** | -0.68 | 0.001 | - | - | - |
| | Ν | 0.029** | -3.84 | -3.08 | 0.06 | 0.008 | - | - | - |
| BSRI | Μ | <u>0.069</u> *** | -6.60*** | -1.87 | 3.88^{***} | 0.113*** | -5.34*** | -4.91*** | 0.21 |
| | F | 0.082^{***} | -6.64*** | -3.05** | 3.04** | 0.109*** | -5.76*** | -4.08*** | 1.04 |
| CABIN | R | 0.110*** | -4.14*** | -6.40*** | -2.91* | 0.078^{***} | -4.87*** | -8.16*** | -4.00*** |
| | Ι | 0.098^{***} | -3.51** | -5.59*** | -3.22** | 0.086^{***} | -4.41*** | -7.75*** | -4.42*** |
| | Α | 0.056^{***} | -3.76*** | -4.63*** | -1.44 | 0.106*** | -4.30*** | -8.00^{***} | -4.14*** |
| | S | 0.092^{***} | -4.05*** | -6.37*** | -3.13** | 0.110*** | -4.70^{***} | -7.60*** | -3.72*** |
| | Ε | 0.081*** | -3.85*** | -5.63*** | -2.44* | 0.117*** | -4.30**** | -8.44*** | -4.31*** |
| | С | 0.048*** | -3.39** | -4.69*** | -1.75 | 0.115*** | -4.95*** | -7.80*** | -3.11** |
| ICB | 0 | 0.025** | -1.83 | -1.49 | 0.22 | 0.073*** | -2.70* | -3.74*** | -1.34 |
| ECR-R | Anx. | 0.236*** | -3.82*** | -8.09*** | -5.33*** | 0.064*** | 0.07 | -2.05 | -2.11 |
| | Avo. | 0.169*** | -3.22** | -7.98*** | -4.61*** | 0.007 | - | - | - |
| MFQ-FF | S. C | 0.063*** | 4.81*** | 4.23*** | -1.09 | 0.007 | - | - | - |
| | Н | 0.067^{***} | 4.95*** | 4.24*** | -1.12 | 0.010 | - | - | - |
| | Ι | 0.071**** | 5.17*** | 4.41*** | -1.26 | 0.007 | - | - | - |
| | R | 0.060 | 4.89 | 4.43*** | -1.06 | 0.005 | - | - | - |
| | S-V | 0.074 | 5.36 | 4.53 | -1.45 | 0.006 | - | - | - |
| | Ε | 0.058 | 5.16 | 4.52 | -1.09 | 0.007 | - | - | - |
| GSE | 0 | 0.074*** | -1.55 | 4.57*** | 6.34*** | 0.039*** | -3.94*** | -3.28** | 0.47 |
| LOT-R | 0 | 0.000 | - | - | - | 0.051*** | -1.91 | -2.83* | -1.37 |
| LMS | R | 0.157*** | -5.85*** | -7.06*** | -2.70^{*} | 0.291*** | -8.11*** | -10.18*** | -4.89*** |
| | М | 0.159*** | -7.23*** | -7.81*** | -2.43* | 0.408^{***} | -8.66*** | -13.20*** | -7.26*** |
| | Ι | 0.196*** | -7.79*** | -8.42*** | -3.30** | 0.449*** | -9.87*** | -14.12*** | -8.18*** |
| EIS | 0 | 0.131*** | -6.93*** | -3.86*** | 2.62^{*} | 0.101*** | -4.84*** | -3.73*** | 0.88 |
| WLEIS | S | 0.080^{***} | -5.28*** | -0.75 | 4.67*** | 0.137*** | -5.33*** | -6.90*** | -2.22 |
| | 0 | 0.021^{*} | -2.95* | 0.14 | 2.87^* | 0.129*** | -5.96*** | -6.87*** | -1.03 |
| | U | 0.073*** | -3.30** | 1.35 | 5.17*** | 0.095*** | -5.06*** | -6.40*** | -1.75 |
| | R | 0.071*** | -3.03** | 2.10 | 5.61*** | 0.147*** | -6.14*** | -7.45*** | -1.47 |
| Empathy | 0 | 0.042*** | -1.88 | -3.65** | -1.99 | 0.004 | - | - | - |
| | | | | | | p | < 0.05 ** p | $p < 0.01^{***}$ | 0 < 0.001 |

Table 14: Result of statistical tests for GPT4o-low and GPT4o-high. Q columns indicate the Q-statistics from the Friedman test (except for GPT4o-low on BSRI Masculine factor, which shows F-statistics from ANOVA, marked with an underline). Also, $\Delta_{i,j}$ columns show the score difference between *i*-th and *j*-th snapshots and corresponding post-hoc test results.

| | Factors | | LLaMA3.1 | 405B-low | | | LLaMA3. | 1 405B-high | |
|---------|---------|---------------------|------------------|------------------|------------------|---------------|------------------|------------------|------------------|
| | | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ | Q | $\Delta_{12,24}$ | $\Delta_{24,36}$ | $\Delta_{12,36}$ |
| BFI | 0 | 0.033** | -1.88 | -2.60^{*} | -1.25 | 0.022^{*} | -1.40 | -2.69* | -1.54 |
| | С | 0.016^{*} | -1.61 | -2.30 | -1.32 | 0.020* | -0.07 | -2.84* | -3.26*** |
| | Ε | 0.012 | - | - | - | 0.019* | -0.48 | -3.05** | -3.14** |
| | Α | 0.025** | -1.98 | -3.06 | -1.89 | 0.034** | -0.54 | -2.56* | -2.60* |
| | Ν | 0.022 | -0.45 | -1.81 | -1.75 | 0.021 | -0.86 | -2.18 | -1.72 |
| EPQ-R | Ε | 0.125^{***}_{***} | 3.07** | -3.57*** | -6.10**** | 0.041*** | -0.84 | -3.91**** | -3.72**** |
| | Р | 0.090*** | 2.37 | -4.42*** | -6.97*** | 0.026** | -0.90 | -4.77*** | -5.04*** |
| | Ν | 0.135 | 2.48 | -5.01 | -6.58 | 0.086 | -1.15 | -5.58 | -5.48 |
| | L | 0.117 | 2.29 | -4.98 | -7.44 | 0.039 | -1.30 | -4.29 | -4.11 |
| DTDD | М | 0.006 | - | - | - | 0.135*** | -4.91*** | -6.67*** | -3.73**** |
| | Р | 0.007 | - | - | - | 0.114*** | -3.82*** | -6.55*** | -4.07**** |
| | Ν | 0.017^{*} | 3.43** | 3.65** | 1.21 | 0.157*** | -1.92 | -7.14*** | -5.55*** |
| BSRI | М | 0.024^{**} | -4.15*** | -1.72 | 2.17 | 0.006 | - | - | - |
| | F | 0.040^{***} | -4.06*** | -2.63* | 1.48 | 0.003 | - | - | - |
| CABIN | R | 0.008 | - | - | - | 0.066*** | -3.47** | -6.57*** | -3.65** |
| | Ι | 0.006 | - | - | - | 0.077^{***} | -3.06** | -4.95*** | -2.33 |
| | A | 0.002 | - | - | - | 0.057^{***} | -3.28** | -4.94*** | -1.92 |
| | S | 0.012 | - | - | - | 0.059^{***} | -4.57*** | -6.36*** | -1.95 |
| | Ε | 0.008 | - | - | - | 0.063*** | -4.54*** | -5.91*** | -1.88 |
| | С | 0.008 | - | - | - | 0.082^{***} | -5.82*** | -5.55*** | -0.51 |
| ICB | 0 | 0.003 | - | - | - | 0.000 | - | - | - |
| ECR-R | Anx. | 0.088^{***} | 1.02 | -6.23*** | -7.88*** | 0.091*** | 2.96^{*} | -3.57** | -7.08*** |
| | Avo. | 0.109*** | -0.12 | -7.35*** | -7.59*** | 0.112*** | 2.05 | -5.12*** | -7.20*** |
| MFQ-FF | S. C | 0.448^{***} | 10.36*** | 11.67*** | 4.49*** | 0.274*** | 3.46** | 9.18*** | 5.82*** |
| | Н | 0.502^{***} | 10.67^{***} | 13.32*** | 5.29^{***} | 0.251*** | 3.45^{**} | 9.57^{***} | 6.32*** |
| | Ι | 0.571^{***} | 11.22^{***} | 13.11*** | 5.14^{***} | 0.357*** | 4.22^{***} | 10.29^{***} | 5.91*** |
| | R | 0.400^{***} | 9.02^{***} | 11.35*** | 4.82^{***} | 0.274^{***} | 4.45^{***} | 9.13*** | 5.77^{***} |
| | S-V | 0.490*** | 11.15^{***} | 12.88*** | 4.55*** | 0.324*** | 4.27^{***} | 10.26^{***} | 6.02^{***} |
| | Ε | 0.440*** | 9.82*** | 11.75*** | 4.63*** | 0.274*** | 3.60** | 9.58*** | 5.10*** |
| GSE | 0 | 0.039*** | -1.81 | 3.54** | 4.84*** | 0.048*** | -1.88 | -4.01*** | -3.42** |
| LOT-R | 0 | 0.025^{**} | 2.14 | 3.48** | 1.82 | 0.024^{**} | -0.21 | -2.32 | -2.47* |
| LMS | R | 0.029** | -2.21 | -3.06** | -1.45 | 0.463*** | -5.34*** | -15.10*** | -12.07*** |
| | М | 0.005 | - | - | - | 0.318*** | -4.01*** | -12.88*** | -9.92*** |
| | Ι | 0.014 | - | - | - | 0.270*** | -3.16** | -11.08*** | -9.35*** |
| EIS | 0 | 0.132*** | -6.89*** | -5.78*** | 1.59 | 0.011 | - | - | - |
| WLEIS | S | 0.056*** | 0.39 | 4.04*** | 3.54** | 0.005 | - | - | - |
| | 0 | 0.025*** | -1.41 | 1.90 | 3.11** | 0.002 | - | - | - |
| | U | 0.043**** | -2.41* | 1.73 | 3.56** | 0.001 | - | - | - |
| | R | 0.018^{*} | -1.05 | 2.09 | 2.78^{*} | 0.000 | - | - | - |
| Empathy | 0 | 0.002 | - | - | - | 0.002 | - | - | - |
| | _ | | | | | * | $p < 0.05^{**}$ | $p < 0.01^{***}$ | p < 0.001 |

Table 15: Result of statistical tests for LLaMA3.1 405B-low and LLaMA3.1 405B-high. Q columns indicate the Q-statistics from the Friedman test. Also, $\Delta_{i,j}$ columns show the score difference between *i*-th and *j*-th snapshots and corresponding post-hoc test results.