Are Economists Always More Introverted? Analyzing Consistency in Persona-Assigned LLMs

Anonymous ACL submission

Abstract

Personalized Large Language Models (LLMs) are increasingly used in diverse applications, where they are assigned a specific persona—such as a happy high school teacher—to guide their responses. While prior research has examined how well LLMs adhere to predefined personas in writing style, a comprehensive analysis of consistency across different personas and task types is lacking. In this paper, we introduce a new standardized framework to analyze consistency in persona-assigned LLMs. We define consistency as the extent to which a model maintains coherent responses when assigned the same persona across different tasks and runs. Our framework evaluates personas across four different categories (happiness, occupation, personality, and political stance) spanning multiple task dimensions (survey writing, essay generation, social media post generation, single turn, and multi-turn conversations). Our findings reveal that consistency is influenced by multiple factors, including the assigned persona, stereotypes, and model design choices. Consistency also varies across tasks, increasing with more structured tasks and additional context. All code is available on GitHub¹.

1 Introduction

011

017

019

021

039

Personalized Large Language Models (LLMs) are increasingly deployed in applications where alignment with specific beliefs and values is essential, such as in high-stakes domains like healthcare and education (Li et al., 2024; Santurkar et al., 2023). While prior research has examined the extent to which LLMs adhere to their assigned personas in terms of writing style (Wang et al., 2024b; Malik et al., 2024), less attention has been given to the consistency of persona adherence across different types of tasks and prompting strategies (Jiang et al., 2024). Moreover, it remains unclear how specifying certain persona attributes affects the consistency of other characteristics. For example, does assigning an "economist" persona to an LLM ensure stable alignment across other characteristics, such as "extroversion"? Additionally, which persona categories lead to the most consistent behavior and does this depend on the task at hand? Addressing these questions is essential for understanding how LLM personas manifest across diverse contexts and for identifying unintended spillover effects-where defining an assigned persona might reinforce unintended other characteristics. Recognizing both the intended and unintended traits associated with a persona is crucial for ensuring reliable and predictable model behavior, especially in high-stake environments.

040

041

042

045

046

047

048

051

052

054

057

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

079

Prior work shows that LLMs can reflect Big Five personality traits in structured tasks, and that larger models tend to do so more consistently than smaller ones (Jiang et al., 2024; Serapio-García et al., 2023). However, persona consistency also extends beyond personality to traits like political orientation and social roles (Röttger et al., 2024; Shu et al., 2024). Yet most evaluations rely on ad hoc methods, such as prompt perturbations or a narrow focus on personality-based personas, resulting in an incomplete picture of consistency. Interestingly, Shu et al. (2024) investigated whether explicitly assigning personas enhances response consistency. They found that while overall consistency decreased with persona assignment, responses became more consistent along dimensions relevant to the assigned persona.

Given the recent calls for application-specific evaluations of LLM behavior (Röttger et al., 2024; Ouyang et al., 2023; Zhao et al.), we propose a new standardized framework for analyzing persona consistency across a broader range of realistic tasks. Moving beyond prompt perturbations and personality-based personas, our approach provides

¹https://anonymous.4open.science/r/persona_ consistency-584C/README.md



Figure 1: The overview of the full methodology. On the left, the persona construction is shown. From these four selected surveys, binary characteristics are selected serving as the base for the different personas. All combinations of these characteristics within a persona category are made, e.g. a character who is introverted, antagonistic, conscientious, neurotic, and open to experience. Surveys are evaluated using their respective scoring key. The other dimensions are evaluated using GPT40 as LLM-as-a-judge.

a systematic analysis of consistency both within and across multiple persona categories and tasks.

More specifically, in this paper, we propose a standardized framework for multifaceted persona consistency analysis. We focus on four persona categories—happiness, occupation, personality, and political stance—selected for their relevance in persona-related literature, their variability in the scale of linguistic expression, and the availability of external survey instruments. We evaluate the consistency of persona-assigned LLMs across multiple evaluation dimensions, including survey answering, social media post generation, essay writing, single-question answering (singlechat), and multi-turn conversations (multichat).

We focus on two primary aspects of consistency: (1) Intra-persona consistency — whether LLMs remain consistent within their assigned persona; and (2) Inter-persona consistency — whether LLMs remain consistent across other persona categories than the one assigned.

We hypothesize that LLMs exhibit personadependent consistency effects: some personas (e.g., tied to specific professions) may induce stronger intra-persona consistency, while others (e.g., based on personality traits) lead to more partial or contextdependent patterns. We also explore potential spillover effects—whether traits associated with one persona category influence outputs in another, possibly due to underlying social stereotypes.

Our findings reveal that specifying a persona

leads to high intra-persona consistency, with some persona categories (e.g., happiness and occupation) being more consistent than others (e.g., political stance). We also uncover spillover effects, where persona assignment reinforces inter-persona consistency driven by stereotypes and model defaults. Finally, we show that consistency is influenced by task dimensions: clearer tasks and additional context length improve consistency. 112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

131

132

133

134

135

136

137

139

2 Framework and Metrics Development

Here, we outline our framework and the evaluation metrics used in our experiments.

2.1 Persona construction

The left part of the consistency framework shown in Figure 1 illustrates the persona construction. We selected the persona categories based on three key criteria: relevance in persona-related literature, variability in the scale of linguistic expression, and the availability of external survey instruments. Specifically, we focused on personality, professions, and political stance (Malik et al., 2024; Jiang et al., 2024; Wang et al., 2024b), which are well-studied categories in the persona-related literature. Additionally, we included a binary persona category—happy or sad—as a useful contrast, since emotional states tend to be more explicitly reflected in language, whereas categories like occupation may manifest more subtly. This range of

- 140 141
- 142
- 143
- 144
- 145 146
- 147
- 148
- 149
- 150 151

- 153 154
- 157

158

- 159

162 163

- 164
- 165
- 167
- 168

169

171

172

173 174

175

176

177 178 179

181 183

184

187

personas allows us to examine varying degrees of consistency. Finally, for each identified persona category, we defined personas based on well-known surveys:

- Happiness: The happiness personas were derived from the Happiness Survey developed by Lyubomirsky and Lepper (1999).
- Political Stance: These personas were based on the Political Compass Test (www. politicalcompass.org/test)
- Occupation: Professional personas were determined using the survey outlined by Holland (1997). More specifically, we chose one occupation per occupation category defined by Holland (1997).
 - Personality: These personas were assigned based on traits from the Big Five Inventory Test (John, 1999).

Based on the outcomes of the surveys, we constructed the different personas by making all possible combinations of the persona characteristics within a persona category, e.g. for the political category we include economically left and socially libertarian; economically right and socially libertarian; economically left and socially authoritarian; economically right and socially authoritarian. Additional information on the selection criteria for the persona categories are provided in Appendix A and all personas are included in Appendix B. This comprehensive approach ensures that our framework captures a wide range of realistic and nuanced persona scenarios.

2.2 Evaluation Dimension selection

The surveys defined in Section 2.1, not only guided the persona construction, but they also serve as one evaluation dimension for assessing the consistency of persona-assigned LLMs. In this evaluation dimension, LLMs are prompted to answer the survey questions individually in separate interactions. Additionally, we identified several other categories to analyze LLM personas: social media post generation, essay writing, single-question answering (singlechat), and multi-turn conversations (multichat). The prompts for these tasks were designed based on the methodologies outlined by Serapio-García et al. (2023) and Jiang et al. (2024). Based on established prompts for social media post generation, we distilled eight separate

open-ended questions as initial prompts for both the singlechat and multichat evaluation dimensions. Multichat, in particular, was specifically designed for this study, building on the same initial prompts as singlechat. After the persona-assigned LLM generated its response, another LLM (LLaMA-3.2-1B) was introduced to engage with the reply. Next, the persona-assigned LLM received the full chat history and was prompted to respond once more. Consistency evaluation was conducted on both responses combined from the persona-assigned LLM. This process is also depicted on the right part of Figure 1. All tasks were carefully selected to align with the call for application-specific evaluations (Röttger et al., 2024; Ouyang et al., 2023; Zhao et al.), ensuring our analysis captures the real-world relevance and practical adaptability of persona-assigned LLMs.

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

226

227

228

229

230

231

232

233

234

235

236

237

2.3 Scoring Key

Additional information on the scoring keys is provided below. A more detailed explanation, including the prompts and specific survey scoring mechanisms, can be found in Appendix C.

Survey Dimension. The survey dimension is evaluated using its respective scoring methodology. To ensure consistency in our analysis, final results are simplified into binary categories for all surveys except the occupational one, where the primary relevant occupational category is selected. For the happiness survey, responses are categorized as happy or sad. In the political compass test, the persona-assigned LLM's outputs are analyzed to determine the corresponding quadrant, with the final outcome identified across both the economic and social axes. In the personality survey, the outputs are assessed within the framework of the Big Five traits and classified into binary categories per trait.

Open Response Dimensions. For the other dimensions, we use an LLM-as-a-judge, GPT4o, to evaluate the final outcomes, determining the persona's alignment with binary characteristics, or for occupations one of the six different classes. The evaluation process is consistent across all dimensions, with the LLM assessing all characteristics across all responses. Detailed information on the used prompt is provided in Appendix C.2. The model also provides a confidence score on a fourpoint Likert scale per choice. A neutral choice was identified when the model's confidence score was

2381 or 2. To validate the reliability of the LLM judg-239ments, we manually annotated 100 examples across240all evaluation categories and found a Cohen's κ of2410.68 with these final LLM results supporting the242reliability of the LLM-generated judgments.

2.4 Consistency scores

244

245

247

248

249

252

253

255

260

261

263

264

267

268

269

270

275

276

277

279

280

281

283

Entropy. To measure consistency, we use Shannon entropy (Shannon, 1948), a metric that captures the uncertainty in the distribution of predicted labels across responses. It applies to both binary and multiclass persona traits and reflects full distributional patterns rather than just majority labels. Crucially, entropy allows us to compare consistency across models and persona categories without relying on arbitrary thresholds.

An entropy score quantifies how focused or scattered the model's responses are. For example, if a model consistently outputs "happy" for a happiness persona across multiple prompts, the entropy is low, indicating high consistency. If the responses are split between "happy" and "sad", the entropy is higher, signaling inconsistency.

We compute entropy for each system prompt s within an evaluation category e, persona category p, and dimension d as:

$$entropy_{s_e,p,d} = -\frac{\sum_{x \in X} P(x) * \log(P(x))}{\log(|X|)}$$
(1)

Here, X is the set of possible characteristics (e.g., for happiness: happy, sad), and P(x) is the proportion of responses labeled with characteristic x. All persona categories are binary except for occupation, which includes six possible labels.

The probability scores per characteristic are calculated using the labels from the LLM-as-a-judge for the different responses. We added the neutral category as a random prediction to every option in the underlying characteristic, because a neutral response indicates a lack of alignment with the intended characteristic. Since consistency requires a persona to manifest in a discernible way, we also treat neutral predictions as inconsistent.

Finally, we compute a single entropy score for each persona and evaluation category by averaging across all system prompts and dimensions:

$$entropy_{p,e} = \frac{1}{|D|} \sum_{d \in D} \frac{1}{|S|} \sum_{s \in S} entropy_{s_e,p,d}$$
(2)

Characteristic-specific consistency. The main disadvantage of the entropy metric is that it does

not show which attribute the LLM consistently outputs. Hence, we also examine the average scores per persona characteristic. For binary characteristics, we use a continuous scale from 0 to 1, where both endpoints represent distinct, persona-aligned responses. A score of 0.5 indicates a lack of alignment with the underlying characteristic, stemming from inconsistency or neutrality. Higher consistency is found when scores are closer to 0 or 1.

285

286

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

For the occupation category, we determine the most frequently assigned occupation category and identify the intensity score as the probability of occurrence. A perfectly consistent model receives an intensity score of 1, while a randomly distributed model is expected to score 1/6.

3 Consistency Analysis

In this section, we present the research questions, the experimental setup, and the results.

3.1 Research Questions

In this study, we examine the consistency of persona-assigned LLMs and spillover effects across different persona categories. Specifically, we address the following research questions:

- 1. **RQ1 Intra-persona consistency:** How does assigning a persona to an LLM result in differences in intra-persona consistency across various persona categories?
- 2. **RQ2 Spillover effects:** Does assigning a persona to an LLM lead to spillover effects in other, unspecified persona categories?
- 3. **RQ3 Cross-dimensional consistency:** How does consistency vary across different response dimensions, particularly with the inclusion of the multichat dimension?

3.2 Experimental set-up

We analyze the consistency over 5 runs across 5 models from 3 different model families: Qwen-2.5 32B, Ministral-8B, Llama-3.2 3B, Llama-3.1 8B, and Llama-3.3 70B. Additional information on the checkpoints used is included in Appendix D We analyze the entropy per model and per evaluation and persona category. To gather insights into the chosen labels per characteristic, we examine the characteristic-specific consistency. Next, we analyze the overall consistency making crossdimension comparisons.

Evaluation Categories		Persona (Categories		Evaluation Categories	Persona Categories				
	Happiness	Occupation	Personality	Political		Happiness	Occupation	Personality	Political	
Happiness	0.18 ± 0.25	0.26 ± 0.41	0.14 ± 0.08	0.21 ± 0.44	Happiness	0.26 ± 0.25	0.27 ± 0.38	0.38 ± 0.16	0.45 ± 0.36	
Occupation	0.59 ± 0.39	0.21 ± 0.21	0.52 ± 0.33	0.51 ± 0.31	Occupation	0.76 ± 0.15	0.38 ± 0.31	0.71 ± 0.15	0.55 ± 0.28	
Personality	0.20 ± 0.14	0.15 ± 0.11	0.25 ± 0.15	0.22 ± 0.11	Personality	0.42 ± 0.15	0.24 ± 0.15	0.38 ± 0.11	0.34 ± 0.08	
Political	0.80 ± 0.45	0.76 ± 0.43	0.75 ± 0.40	0.41 ± 0.46	Political	0.86 ± 0.19	0.86 ± 0.16	0.51 ± 0.35		
(a) Ent	ropy score	s for Qwer	n-2.5 32B.		(b) En	(b) Entropy scores for Ministral-8B.				
Evaluation Categories		Persona (ategories		Evaluation Categories		Persona (ategories		
Evaluation Categories		Persona (Categories		Evaluation Categories		Persona (Categories		
Evaluation Categories	Happiness	Persona Occupation	Categories Personality	Political	Evaluation Categories	Happiness	Persona C	Categories Personality	Political	
Evaluation Categories Happiness	Happiness 0.00 ± 0.00	Persona C Occupation 0.30 ± 0.25	Categories Personality 0.54 ± 0.12	Political 0.64 ± 0.17	Evaluation Categories Happiness	Happiness 0.07 ± 0.16	Persona O Occupation 0.21 ± 0.14	Categories Personality 0.43 ± 0.09	Political 0.39 ± 0.14	
Evaluation Categories Happiness Occupation	Happiness 0.00 ± 0.00 0.73 ± 0.16	Persona 0 Occupation 0.30 ± 0.25 0.42 ± 0.24	Categories Personality 0.54 ± 0.12 0.66 ± 0.21	Political $\begin{array}{c} 0.64 \pm 0.17\\ 0.63 \pm 0.23 \end{array}$	Evaluation Categories Happiness Occupation	Happiness 0.07 ± 0.16 0.56 ± 0.30	Persona C Occupation 0.21 ± 0.14 0.23 ± 0.19	Categories Personality 0.43 ± 0.09 0.66 ± 0.24	Political 0.39 ± 0.14 0.59 ± 0.21	
Evaluation Categories Happiness Occupation Personality	Happiness 0.00 ± 0.00 0.73 ± 0.16 0.31 ± 0.19	Personal Occupation 0.30 ± 0.25 0.42 ± 0.24 0.31 ± 0.12	Categories Personality 0.54 ± 0.12 0.66 ± 0.21 0.42 ± 0.11	Political 0.64 ± 0.17 0.63 ± 0.23 0.43 ± 0.16	Evaluation Categories Happiness Occupation Personality	Happiness 0.07 ± 0.16 0.56 ± 0.30 0.30 ± 0.15	Persona C Occupation 0.21 ± 0.14 0.23 ± 0.19 0.22 ± 0.04	CategoriesPersonality 0.43 ± 0.09 0.66 ± 0.24 0.35 ± 0.13	Political 0.39 ± 0.14 0.59 ± 0.21 0.41 ± 0.12	
Evaluation Categories Happiness Occupation Personality Political	Happiness 0.00 ± 0.00 0.73 ± 0.16 0.31 ± 0.19 0.84 ± 0.23	$\begin{array}{c} Persona \ 0\\ \hline 0.30 \pm 0.25\\ 0.42 \pm 0.24\\ 0.31 \pm 0.12\\ 0.81 \pm 0.25 \end{array}$	Categories Personality 0.54 ± 0.12 0.66 ± 0.21 0.42 ± 0.11 0.89 ± 0.13	$\begin{tabular}{l} \hline Political \\ 0.64 \pm 0.17 \\ 0.63 \pm 0.23 \\ 0.43 \pm 0.16 \\ 0.70 \pm 0.18 \end{tabular}$	Evaluation Categories Happiness Occupation Personality Political	Happiness 0.07 ± 0.16 0.56 ± 0.30 0.30 ± 0.15 0.83 ± 0.33	$\begin{array}{c} Persona \ 0\\ \hline 0.21 \pm 0.14 \\ 0.23 \pm 0.19 \\ 0.22 \pm 0.04 \\ 0.75 \pm 0.38 \end{array}$	Categories Personality 0.43 ± 0.09 0.66 ± 0.24 0.35 ± 0.13 0.79 ± 0.30	Political 0.39 ± 0.14 0.59 ± 0.21 0.41 ± 0.12 0.36 ± 0.31	

Evaluation Categories	Persona Categories								
	Happiness	Occupation	Personality	Political					
Happiness	0.07 ± 0.15	0.05 ± 0.09	0.30 ± 0.19	0.15 ± 0.10					
Occupation	0.62 ± 0.38	0.23 ± 0.20	0.56 ± 0.34	0.49 ± 0.34					
Personality	0.23 ± 0.16	0.16 ± 0.09	0.27 ± 0.18	0.26 ± 0.13					
Political	0.79 ± 0.44	0.74 ± 0.43	0.71 ± 0.41	0.40 ± 0.31					

(e) Entropy scores for Llama-3.3 70B.

Table 1: The tables show large entropy differences between the models, indicating differences in consistency levels. Both strong within-category consistency (diagonal) and occasional spill-over consistency (off-diagonal) are found, where lower is more consistent. Scores <0.25 are colored green, 0.25<0.5 in orange, and >0.5 in red. The columns represent the different assigned persona categories. The rows represent the evaluation categories. The standard deviation is computed over the entropy scores over different dimensions per evaluation-persona category pair.

3.3 Results

Intra-persona consistency is high within each category, but notable differences emerge across categories (RQ1). As shown in Table 1, the diagonal values indicate relatively high consistency, meaning that when a persona is assigned, the model tends to generate responses that consistently express that specific persona across different output formats and prompts. However, the degree of consistency varies across persona categories. While happiness and occupation personas are more consistently expressed, personality and political personas exhibit lower intra-persona consistency. For the political category, we observe a high standard deviation, indicating substantial variability in consistency across dimensions. This is largely due to certain tasks, such as singlechat, where expressing a consistent political stance is more challenging in general. Manual analysis also confirms this finding, highlighting the difficulty of a model to express the political opinion when being asked certain questions. Here, the inconsistency thus stems from a lack of expression of the underlying persona. Similarly, personality-based personas show lower intra-persona consistency, which we investigate further in Figure 2. This figure shows the results for Qwen-2.5 32B (other models are shown in Appendix E). We chose Qwen as it provides representative results for all models with an average correlation of 0.70 with the other models. Examining Figure 2, we find that the LLM generally follows our instructions, e.g., the happy persona is more happy (1), while the sad persona is more sad (0). Likewise assigned occupations are clearly reflected in the output. However, certain personality traits-such as low conscientiousness, antagonism, neuroticism, and, in some models, low openness to experience—are more difficult for the LLM to express consistently, resulting in greater variability in responses within that personality category. Similarly, Salecha et al. (2024) show that models skew responses to socially desirable answers for the Big Five Personality test when they infer that they are evaluated. We demonstrate that social desirability bias also appears in other evaluation dimensions. This social desirability tendency explains the model's difficulty to adhere to the unconscientious and antagonistic personas.

359

360

361

362

363

364

365

366

367

368

369

371

372

373

374

376

377

378

379

380

381

382

384

385

Spill-over effects vary across evaluation categories (RQ2). The off-diagonal elements in the subtables of Table 1 reveal that most spill-over consistency effects occur across two evaluation categories: happiness and personality. For example, the off-diagonal entropy scores for Llama-3.3 70B are lower for the happiness and personality categories compared to occupation and political stance.

353

354

357

331



(a) Heatmap providing the characteristic-specific consistency for all evaluation categories except occupation. A score of 1 favors the category name, 0 favors its opposite (e.g., agreeableness vs. antagonistic), and 0.5 indicates inconsistency.



(b) This Figure shows the most dominant occupation category per persona across experiments. Color intensity represents the consistency score per label.

Figure 2: Qwen-2.5 32B generally follows the instructions and shows spill-over effects, i.e. stereotypes and default values. Columns denote personas, and rows indicate evaluation categories. Multi-component personas (e.g., political stance) are grouped per component and averaged over all personas containing that component. The other models are shown in Appendix E.

This suggests that the model is more consistent across these two categories when they are not the assigned persona. Interestingly, this pattern differs from the intra-persona results, where happiness and occupation showed the strongest consistency when they were the assigned persona. The occupation and political categories are harder for the models to express, as not adding those personas results in responses without any occupational information or political stance. Manual analysis reveals how these personas are less frequently and explicitly expressed, especially in conversational settings. For instance, political beliefs rarely surface in responses to questions like "What are your music preferences?". Similarly, occupation-related information rarely appears unless explicitly prompted, though it may occasionally surface in essay-style or social media posts. The other two categories show spill-over effects. Additional manual analysis re-

387

391

400

401

402

403

404

veals that happiness and personality directly influence linguistic style—models default to a positive tone unless instructed otherwise, making happiness more overt. Similarly, personality traits, like extroversion, shape response style, amplifying spillover effects. 405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

Spill-over effects are due to two main factors: stereotypical associations with the assigned persona and default personas when a characteristic is not explicitly assigned (RQ2). Figure 2 shows how a sad persona is portrayed as more introverted, less conscientious, and less open than a happy persona. The sad persona is more likely to have an artistic occupation, while the happy persona is most likely to have a social occupation. The economically right-winged and socially authoritarian personas are both less agreeable than their counterparts. All occupation personas are presented as extroverts, except the economist, who is more introverted. These observations reinforce prior findings that persona-assigned LLMs are susceptible to stereotypes (Gupta et al., 2023). We show that these stereotypes also appear across general text-generation tasks. Furthermore, the model tends to answer in a happy, conscientious, agreeable, and open manner unless otherwise instructed or influenced by a stereotype. This is reflected in the heatmaps: the rows corresponding to the happiness, conscientiousness, and agreeableness evaluation categories generally lean toward a value of 1, indicating strong alignment. In contrast, neuroticism tends to lean toward 0, suggesting lower identification with that trait. This again illustrates the social desirability bias in the models Salecha et al. (2024). Additionally, we show how personas can partially counteract this bias, e.g. the high consistency scores for neurotic and sad personas. However, this effect is not universal, as evidenced by the lower intra-persona consistency for antagonistic and unconscientious personas. Furthermore, most models tend to be slightly economically left and socially libertarian, though this varies by model and sometimes leans toward inconsistency. These default personas reveal the LLM's ideological stances, shaped by training data and choices from model developers, called design choices (Buyl et al., 2025; Cambo and Gergle, 2022).

Consistency is higher in more structured tasks like survey answering. For open-ended question answering tasks, providing additional context through multi-turn interactions (multichat)



Figure 3: The average intra-persona and inter-persona entropy across all models per dimension reveals large differences in entropy between different dimensions. Error bars represent 95% confidence intervals estimated via bootstrapping, using nonparametric resampling to approximate the uncertainty around the mean.

improves consistency (RQ3). Figure 3 shows 456 457 that the survey, essay, and social media dimensions have the lowest intra- and inter-persona entropy, in-458 dicating the highest consistency across these tasks. 459 However, consistency scores vary notably across 460 dimensions, with mostly low correlations between 461 462 them. This highlights the importance of considering all three dimensions to fully capture model 463 behavior. Of these, only the inter-persona consis-464 tency between essay and social media is strongly 465 correlated. Their low entropy scores could be at-466 tributed to the clarity of the task, where models 467 can easily express their assigned persona. As tasks 468 become less straightforward—such as answering 469 open-ended questions about music preferences (sin-470 glechat and multichat)-models generally show a 471 decrease in both intra- and inter-persona consis-472 tency. Moreover, the difference between intra- and 473 inter-persona entropy becomes smaller. These re-474 sults suggest that as tasks require less explicit per-475 sona expression, models may struggle to express 476 distinct persona characteristics. Depending on the 477 application, this variability may be advantageous, 478 i.e. in creative or open-ended generation tasks 479 where diversity is desirable. However, in more 480 controlled settings that require reliable persona ad-481 herence, such unpredictability can be a drawback, 482 as the model's outputs may deviate from the in-483 tended persona or produce inconsistent behavior. 484 485 Interestingly, the complexity of multichat scenarios compared to singlechat does not appear to hin-486 der consistency. Contrarily, consistency increases 487 slightly as follow-up responses allow models to 488 provide more information, expressing the assigned 489

persona more clearly. Nevertheless, consistency scores between singlechat and multichat are highly correlated. In addition to the main results, supplementary statistical analyses revealed no significant difference between inter- and intra-persona consistency in the singlechat and multichat dimensions. However, we observed statistically significant differences between the evaluation dimensions of survey, essay, and social media on the one hand, and singlechat (for both intra- and inter-persona consistency) as well as multichat (for intra-persona consistency) on the other. Concerning inter-persona consistency specifically, only the social media post generation task significantly outperforms the interpersona consistency observed in the multichat evaluation dimension. Detailed statistical analysis is reported in Appendix G.

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

4 Discussion

Our framework offers a multi-dimensional perspective on the consistency of persona-assigned LLMs. Our results show that the balance between intrapersona and inter-persona consistency varies depending on the evaluation dimension. Personabased models are more consistent for attributes explicitly defined in their persona than for unspecified persona attributes, indicating a weaker spillover effect. However, in singlechat and multichat settings, we find no statistical difference between inter- and intra-persona consistency, implying that adding a persona does not significantly increase consistency compared to non-assigned attributes in more realistic tasks. This finding urges caution when deploying persona-assigned LLMs in practical applications, as persona adherence may be less consistent in conversational contexts than in more structured tasks like surveys or essays. Finally, we show that longer context, i.e. comparing multichat to singlechat, allows models to express their persona more clearly leading to higher consistency.

We identify three main factors influencing consistency across different characteristics.

The assigned persona: Models generally adhere well to persona-specific instructions, with the degree of adherence varying across categories, models and dimensions. Our paper offers a valuable comparison on consistency levels across different persona categories, highlighting how some categories (e.g. happiness and occupation) are easier to consistently express than others (e.g. personality and political stance). The evaluation dimension

631

632

633

634

635

636

637

638

also highly influences consistency. We show how 540 more structured tasks and longer sequence lengths 541 in the response result in higher consistency, espe-542 cially for the identified harder persona categories, such as personality and political stance.

545 Stereotypical associations with the assigned persona: Stereotypical associations play a significant role in inter-persona consistency. Characteristics that align with stereotypical traits of a given persona often result in higher consistency scores. For example, persona-assigned LLMs instructed to be "happy" consistently score higher on extroversion. These tendencies highlight the influence of societal stereotypes embedded within the models. Literature confirms this influence of stereotypes in persona-assigned LLMs (Gupta et al., 2023).

547

551

554

556

557

560

563

567

569

571

572

573

575

577

584

585

588

Default persona: When a specific characteristic is not defined and a persona lacks a stereotypical association, the model often reverts to a consistent, pre-defined default persona. Models exhibit social desirability bias as was found by Salecha et al. (2024). Our findings also reveal default political personas in these models, likely reflecting design choices by developers. As shown by Buyl et al. (2025), ideological stances can be embedded in models, a phenomenon tied to model positionality-the social and cultural perspective developers align the model with (Cambo and Gergle, 2022). These choices occur throughout the development process and extend beyond training data alone (Buyl et al., 2025).

5 **Related work**

Persona-assigned LLMs. Personas can guide LLMs to generate responses that align with specific values and beliefs (Li et al., 2024; Santurkar et al., 2023). However, they can also expose stereotypes embedded in the model (Park et al., 2024; Gupta et al., 2023), raising concerns about bias and unintended implications. Persona adherence is usually evaluated using self-report scales, but Wang et al. (2024b) use interview-based testing to capture actual model behavior, showing the need for application-specific evaluations. Malik et al. (2024) examine how different personas from various sociodemographic groups influence writing styles. Apart from inference-time persona assignment, it is also possible to further fine-tune the model. For example, Shao et al. (2023) train LLMs to adopt specific personas using three key components: a

profile (a detailed persona description), a scene (a situational context), and interactions. Wang et al. (2024a) suggest including dialogues for persona assignment of LLMs. Our framework can evaluate these realistic personas or finetuned models that do not perfectly fit our predefined categories, as shown in Appendix H.

LLM consistency. LLMs are self-inconsistent when prompted with ambiguous entities (Sedova et al., 2024). Röttger et al. (2024) show how models do not answer consistently when paraphrasing prompts from the political compass test. Shu et al. (2024) show how LLMs are inconsistent over different prompt perturbations. When analyzing the effect of adding a persona when measuring model consistency, overall assigning a persona does not help consistency. Nevertheless, consistency improves along the axes relevant to the persona. Recently, Lee et al. (2024) introduced a multiple-choice benchmark dataset to assess consistency in LLM outputs with respect to personality. While their analysis focuses on consistency within model responses using their dataset, it is limited to multiple-choice scenarios-aligning with our survey evaluation dimension—whereas we have shown that consistency can vary significantly across different evaluation dimensions. Finally, Jiang et al. (2024) evaluate whether persona-assigned LLMs consistently follow personality traits from the Big Five personality test for two evaluation dimensions: survey and essay.

6 Conclusion

This paper introduces a multi-dimensional framework for analyzing consistency in persona-assigned LLMs. Our framework encompasses a diverse set of commonly used persona categories, including personality, occupation, political stance, and happiness. It also incorporates application-specific evaluation dimensions, such as survey answering, essay writing, social media post generation, singlequestion answering, and multi-turn conversations. We demonstrate the efficacy of our framework through a comprehensive evaluation of intra- and inter-persona consistency across personas derived from the defined persona categories. Additionally, we compare consistency scores across evaluation dimensions. Our analysis reveals three key factors influencing consistency in LLMs: (1) the assigned persona; (2) stereotypes associated with the assigned persona; and (3) model's default personas.

7 Limitations

639

640

641

643

646

647

667

670

672

673

675

677

684

We have used an LLM-as-a-judge for the annotations of our results. However, these models are very sensitive towards several different types of biases. It is known that LLMs can be subject to order bias (Li et al., 2025). By adding confidence scores, we have mitigated this bias partly. We have tested it on a subsample of our dataset and found that there was indeed order bias, however, this mainly occurred when there was a low confidence in the given answer. The Cohen's kappa of a manually validated sample and the sample used in our paper was 65.42%, for the sample where orders were reversed, the Cohen's kappa was 65.49%. We thus assume this did not influence our results that much. We also only used one LLM-as-a-judge for our analyses. We checked for other LLMs on a subsample and they performed similarly. Here we found a Cohen's kappa of 69.64% for Sonnet on a sample of 100 manual validations. Moreover, to avoid self-preference bias within LLMs (Li et al., 2025), we used different LLMs than the ones that we used for the first answer generation. Moreover, further analysis, including additional LLMs and focusing on how architectural and training differences impact consistency in LLMs would be a valuable direction for future work. Finally, future work could investigate the impact of post-training alignment on role-playing capabilities. In particular, comparing instruction-tuned models with their base counterparts may offer deeper insights into how alignment influences persona consistency.

8 Ethical considerations

We should be aware when using LLM-as-a-judge that there exists demographic bias towards certain groups, especially in subjective tasks as is shown in (Alipour et al., 2024). Furthermore, this paper highlights how LLMs have been trained with certain design choices. So when a value is not explicitly described, they tend to go to a certain default persona. It is important to keep in mind that this will differ across different LLMs. Additionally, the stereotypes learned by the model and also consistently expressed are thus also very model-specific. Moreover, it is important to note that consistency is not the same as ethical correctness. Therefore, there is still a need for responsible model deployment even though models might already provide rather consistent answers. Finally, people should be aware when adding personas to LLMs that certain stereotypes

might be inherently present in these models, further reinforcing certain stereotypes.

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

References

- Shayan Alipour, Indira Sen, Mattia Samory, and Tanushree Mitra. 2024. Robustness and confounders in the demographic alignment of llms with human perceptions of offensiveness. *arXiv preprint arXiv:2411.08977*.
- Maarten Buyl, Alexander Rogiers, Sander Noels, Guillaume Bied, Iris Dominguez-Catena, Edith Heiter, Iman Johary, Alexandru-Cristian Mara, Raphaël Romero, Jefrey Lijffijt, and Tijl De Bie. 2025. Large language models reflect the ideology of their creators. *Preprint*, arXiv:2410.18417.
- Scott Allen Cambo and Darren Gergle. 2022. Model positionality and computational reflexivity: Promoting reflexivity in data science. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–19.
- Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,

Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal 746 Lakhotia, Lauren Rantala-Yeary, Laurens van der 747 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, 749 Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, 754 Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Celebi, Patrick 757 Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, 761 Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten 771 Sootla, Stephane Collot, Suchin Gururangan, Syd-773 ney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias 774 775 Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal 776 Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh 777 Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petro-778 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-779 ney Meers, Xavier Martinet, Xiaodong Wang, Xi-781 aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei 790 Baevski, Allie Feinstein, Amanda Kallet, Amit San-791 gani, Amos Teo, Anam Yunus, Andrei Lupu, An-792 dres Alvarado, Andrew Caples, Andrew Gu, Andrew 793 Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-794 dani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, 796 Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-797 dan, Beau James, Ben Maurer, Benjamin Leonhardi, 798 Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-799 cock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, 803 Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, 808 Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,

Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The Ilama 3 herd of models. *Preprint*, arXiv:2407.21783.

874

875

893

894

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

921

922

923

924

925

926

927

928

929 930

- Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2023. Bias runs deep: Implicit reasoning biases in persona-assigned llms. In *The Twelfth International Conference on Learning Representations*.
 - John L. Holland. 1997. *Making vocational choices: a theory of vocational personalities and work environments.*, third edition edition. PAR, Lutz.
 - Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2024. PersonaLLM: Investigating the ability of large language models to express personality traits. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3605–3627. Association for Computational Linguistics.
 - OP John. 1999. The big-five trait taxonomy: History, measurement, and theoretical perspectives. *Handbook of Personality: Theory and Research/Guilford.*
 - Seungbeen Lee, Seungwon Lim, Seungju Han, Giyeong Oh, Hyungjoo Chae, Jiwan Chung, Minju Kim, Beong-woo Kwak, Yeonsoo Lee, Dongha Lee, et al. 2024. Do llms have distinct and consistent personality? trait: Personality testset designed for llms with psychometrics. *CoRR*.
 - Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. 2025. From generation to judgment: Opportunities and challenges of llm-as-a-judge. *Preprint*, arXiv:2411.16594.
 - Junyi Li, Charith Peris, Ninareh Mehrabi, Palash Goyal, Kai-Wei Chang, Aram Galstyan, Richard Zemel, and Rahul Gupta. 2024. The steerability of large language models toward data-driven personas. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 7290–7305. Association for Computational Linguistics.
 - Andy Liu, Mona Diab, and Daniel Fried. 2024. Evaluating large language model biases in persona-steered generation. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 9832–9850. Association for Computational Linguistics.

Sonja Lyubomirsky and Heidi S Lepper. 1999. A measure of subjective happiness: Preliminary reliability and construct validation. *Social indicators research*, 46:137–155. 931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

- Manuj Malik, Jing Jiang, and Kian Ming A. Chai. 2024. An empirical analysis of the writing styles of personaassigned LLMs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19369–19388, Miami, Florida, USA. Association for Computational Linguistics.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,

Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.

993

994

997

999

1002

1003

1004

1005

1008

1011

1013

1014

1015

1016

1018

1020

1021

1023

1024

1025

1026

1027

1028

1029

1034

1035 1036

1037

1038

1040

1043 1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

- Siru Ouyang, Shuohang Wang, Yang Liu, Ming Zhong, Yizhu Jiao, Dan Iter, Reid Pryzant, Chenguang Zhu, Heng Ji, and Jiawei Han. 2023. The shifted and the overlooked: A task-oriented investigation of user-GPT interactions. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2375–2393. Association for Computational Linguistics.
- Joon Sung Park, Carolyn Q Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S Bernstein.
 2024. Generative agent simulations of 1,000 people. arXiv preprint arXiv:2411.10109.
- Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Kirk, Hinrich Schuetze, and Dirk Hovy. 2024. Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15295–15311. Association for Computational Linguistics.
- Aadesh Salecha, Molly E Ireland, Shashanka Subrahmanya, João Sedoc, Lyle H Ungar, and Johannes C Eichstaedt. 2024. Large language models display

human-like social desirability biases in big five personality surveys. *PNAS Nexus*, 3(12):pgae533.

1054

1055

1056

1057

1058

1059

1060

1061

1063

1064

1065

1066

1067

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1082

1083

1084

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

- Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023.
 Whose opinions do language models reflect? In *International Conference on Machine Learning*, pages 29971–30004. PMLR.
- Anastasiia Sedova, Robert Litschko, Diego Frassinelli, Benjamin Roth, and Barbara Plank. 2024. To know or not to know? analyzing self-consistency of large language models under ambiguity. In *Findings of the Association for Computational Linguistics: EMNLP* 2024, pages 17203–17217, Miami, Florida, USA. Association for Computational Linguistics.
- Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality traits in large language models. *Preprint*, arXiv:2307.00184.
- Claude Elwood Shannon. 1948. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-LLM: A trainable agent for roleplaying. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13153–13187. Association for Computational Linguistics.
- Bangzhao Shu, Lechen Zhang, Minje Choi, Lavinia Dunagan, Lajanugen Logeswaran, Moontae Lee, Dallas Card, and David Jurgens. 2024. You don't need a personality test to know these models are unreliable: Assessing the reliability of large language models on psychometric instruments. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 5263–5281. Association for Computational Linguistics.
- Qwen Team. 2024. Qwen2.5: A party of foundation models.
- Noah Wang, Z.y. Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. 2024a. RoleLLM: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 14743–14777. Association for Computational Linguistics.
- Xintao Wang, Yunze Xiao, Jen-tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, Jiangjie Chen, Cheng Li, and Yanghua Xiao. 2024b. InCharacter: Evaluating personality fidelity in role-playing agents through psychological interviews. In *Proceedings of the 62nd*

- 1111 1112 1113 1114 1115 1116 1117 1118 1119
- 1120 1121
- 1122

stance.

Α

- 1123
- 1124 1125
- 1126
- 1127 1128
- 1129 1130
- 1131 1132
- 1133
- 1134

1135 1136

- 1137 1138
- - -

1139 1140

1141 1142

1143 1144

1145 1146

1147

1148

1149

1150 1151

1

We acknowledge that these categories do not rep-1152 resent an exhaustive set of personas. They were 1153 chosen based on the previously mentioned criteria, 1154 1155 balancing theoretical grounding in persona literature and social science research with linguistic 1156 variation. RQ5 further examines the real-world ap-1157 plicability by incorporating diverse personas from 1158 the PersonaHub. 1159

trait, as a baseline for comparison.

Annual Meeting of the Association for Computational

Linguistics (Volume 1: Long Papers), pages 1840–

1873. Association for Computational Linguistics.

Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie,

Mingqian Zheng, Jiaxin Pei, and David Jurgens. 2023.

Is" a helpful assistant" the best role for large language

models? a systematic evaluation of social roles in

system prompts. arXiv preprint arXiv:2311.10054,

As said, we chose the persona categories based on

three main criteria. Below we will extend upon

Relevance in existing research: We analyzed

prior work on persona-related studies and identified

commonly used categories used in the NLP/ACL

community. Many studies focus on personality

traits, often framed through the Big Five model

(e.g., Jiang et al. (2024); Serapio-García et al.

(2023)). Similarly, political stance has been ex-

tensively studied (e.g., Liu et al. (2024)), as well

as social roles and occupations (e.g., Zheng et al.

(2023)). Moreover, Röttger et al. (2024) exam-

ine consistency of LLMs regarding their political

Variation in linguistic expression: We include

persona categories that differ in how explicitly they

manifest in language. For instance, occupations

and political stances may be challenging to express

consistently, while personality traits might influ-

ence linguistic style more subtly. To complement

these, we included happiness, a simpler, binary

Grounding in social science literature: To

ensure robustness, we selected persona categories

with established survey instruments, allowing for

structured and reproducible persona assignment.

Persona Categories

each of these criteria.

Yejin Choi, and Yuntian Deng. Wildchat: 1m chatgpt

interaction logs in the wild. In *The Twelfth Interna*tional Conference on Learning Representations.

B Persona prompts

All personas are added to the system prompt using1161the following template: You are a character who is1162[persona]. The different personas are listed below1163per persona category.1164

1160

1165

1166

1167

1168

1169

1170

1171

Happiness: happy, sad

Political: economically left and socially libertarian; economically left and socially authoritarian; economically right and socially libertarian; economically right and socially authoritarian

Occupations: a pilot; an economist; an actor; a nurse; a CEO; a sales assistant

Personality: introverted, antagonistic, unconsci-1172 entious, neurotic, and open to experience; intro-1173 verted, antagonistic, unconscientious, neurotic, and 1174 closed to experience; introverted, antagonistic, un-1175 conscientious, emotionally stable, and open to 1176 experience; introverted, antagonistic, unconscien-1177 tious, emotionally stable, and closed to experience; 1178 introverted, antagonistic, conscientious, neurotic, 1179 and open to experience; introverted, antagonistic, 1180 conscientious, neurotic, and closed to experience; 1181 introverted, antagonistic, conscientious, emotion-1182 ally stable, and open to experience; introverted, an-1183 tagonistic, conscientious, emotionally stable, and 1184 closed to experience; introverted, agreeable, uncon-1185 scientious, neurotic, and open to experience; intro-1186 verted, agreeable, unconscientious, neurotic, and 1187 closed to experience; introverted, agreeable, un-1188 conscientious, emotionally stable, and open to ex-1189 perience; introverted, agreeable, unconscientious, 1190 emotionally stable, and closed to experience; in-1191 troverted, agreeable, conscientious, neurotic, and 1192 open to experience; introverted, agreeable, consci-1193 entious, neurotic, and closed to experience; intro-1194 verted, agreeable, conscientious, emotionally sta-1195 ble, and open to experience; introverted, agreeable, 1196 conscientious, emotionally stable, and closed to 1197 experience; extroverted, antagonistic, unconscien-1198 tious, neurotic, and open to experience; extroverted, 1199 antagonistic, unconscientious, neurotic, and closed 1200 to experience; extroverted, antagonistic, unconsci-1201 entious, emotionally stable, and open to experi-1202 ence; extroverted, antagonistic, unconscientious, 1203 emotionally stable, and closed to experience; extro-1204 verted, antagonistic, conscientious, neurotic, and 1205 open to experience; extroverted, antagonistic, con-1206 scientious, neurotic, and closed to experience; extroverted, antagonistic, conscientious, emotionally 1208

stable, and open to experience; extroverted, an-1209 tagonistic, conscientious, emotionally stable, and 1210 closed to experience; extroverted, agreeable, uncon-1211 scientious, neurotic, and open to experience; extro-1212 verted, agreeable, unconscientious, neurotic, and closed to experience; extroverted, agreeable, un-1214 conscientious, emotionally stable, and open to ex-1215 perience; extroverted, agreeable, unconscientious, 1216 emotionally stable, and closed to experience; ex-1217 troverted, agreeable, conscientious, neurotic, and 1218 open to experience; extroverted, agreeable, con-1219 scientious, neurotic, and closed to experience; ex-1220 troverted, agreeable, conscientious, emotionally 1221 stable, and open to experience; extroverted, agree-1222 able, conscientious, emotionally stable, and closed 1223 to experience 1224

C Evaluation

1225

1226

1227

1228

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

In this section, more information is given about the evaluation methods used, including the different survey scoring keys and the prompt for the LLMas-a-judge.

C.1 Survey

Happiness: To determine the final label, we first computed a total score based on the generated responses using the Likert scale. We then applied a thresholding approach, assigning a binary label depending on whether the score fell above or below the midpoint.

Political: To determine the final label, we first computed a total score per characteristic (economic and social stance) based on the generated responses and the survey scoring mechanism. We then applied a thresholding approach, assigning a binary label depending on whether the score fell above or below the midpoint.

1244Occupations: We calculated a score for every oc-1245cupation category based on the generated responses1246and the survey scoring mechanism. Next, we identi-1247fied the category with the highest score as the main1248occupational category.

1249**Personality:** Each statement in the survey cor-1250responds to a particular binary characteristic. To1251determine the final trait labels, we first computed a1252total score for each trait based on the generated re-1253sponses. We then applied a thresholding approach,1254assigning a binary label depending on whether the1255score fell above or below the midpoint.

C.2 LLM-as-a-judge

The system prompt given to the LLM-as-a-judge is the following: You will help annotating a dataset. Answer the questions as asked, do not provide extra explanations, only choose one of the provided options. 1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

The following user prompt is used for evaluating the **essay, social media, and singlechat** dimensions.

Below you can find a question and the answer is written by a person. [instruction] [LLM response] Is the person who wrote this [happy or sad, extrovert or introvert, agreeable or antagonistic, conscientious or unconscientious, neurotic or emotionally stable, open to experiences or closed to experiences, economically left or economically right, socially libertarian or socially authoritarian, having a realistic occupation (i.e. pilot) or an investigative occupation (i.e. economist) or an artistic occupation (i.e. actor) or a social occupation (i.e. nurse), or an enterprising occupation (i.e. CEO), or a conventional occupation (i.e. sales assistant)]

You **must** between one of the options. First output this category, then generate a confidence score, where 1 is Very unconfident, 2 is unconfident, 3 is confident, 4 is very confident. Use the following structured output and fill in [category] and [confidence_score]: choice:[category],confidence:[confidence_score]

For the **multichat** dimension, the following 1289 prompt is used: Below you can find the conver-1290 sation between an assistant and person 1291 assistant: [instruction] person: [response] 1293 assistant: [response assistant] 1294 person: [final response] 1295 Is the person who wrote this [happy or sad, extro-1296 vert or introvert, agreeable or antagonistic, consci-1297 entious or unconscientious, neurotic or emotionally 1298 stable, open to experiences or closed to experiences, 1299 economically left or economically right, socially 1300 libertarian or socially authoritarian, having a re-1301 alistic occupation (i.e. pilot) or an investigative 1302 occupation (i.e. economist) or an artistic occupa-1303 tion (i.e. actor) or a social occupation (i.e. nurse), 1304 or an enterprising occupation (i.e. CEO), or a con-1305 ventional occupation (i.e. sales assistant)] 1306 1307You **must** between one of the options. First1308output this category, then generate a confi-1309dence score, where 1 is Very unconfident, 21310is unconfident, 3 is confident, 4 is very con-1311fident. Use the following structured output1312and fill in [category] and [confidence_score]:1313choice:[category],confidence:[confidence_score]

D Model Checkpoints

1314

1327

1328

1331

1332

1334

1335

1336

1337

1338

1339

1340

1342

1343

1344

1346

1347

1348

1349

For the different experiments we used the follow-1315 ing model checkpoints: meta-llama/Llama-3.2-3B-1316 Instruct, meta-llama/Llama-3.1-8B-Instruct, meta-1317 llama/Llama-3.3-70B-Instruct (Grattafiori et al., 1318 2024), mistralai/Ministral-8B-Instruct-2410², and 1319 Qwen/Qwen2.5-32B-Instruct-GPTQ-Int8 (Team, 1320 2024). For the LLM-evaluation, we used gpt-4o-2024-08-06 (OpenAI et al., 2024). We ran our 1322 experiments on H100 GPUs. All models were used 1323 consistent with their intended use and in line with 1324 their provided licenses. The temperature for all 1325 experiments was set at 0.7. 1326

E Characteristic-specific consistency

Figures 4, 5, 6, and 7 provide insights into characteristic-specific consistency of Llama 3B, Llama 8B, Llama 70B, and Ministral respectively.

F Model Comparisons

In this section we focus on the following research question: Do consistency patterns differ across model families and/or within a single model family?

Consistency varies across model families and within a model family it increases with model size (RQ4). Figure 8 illustrates how consistency varies across model families. For example, we see that Ministral-8B has lower overall consistency than Llama-3.1-8B despite similar model size. Additionally, our results show that within a model family, larger models tend to be more consistent than smaller models. This is shown by the three Llama models in the figure. This result aligns with the finding from Serapio-García et al. (2023), showing higher reliability and validity of synthetic LLM personality for larger and instruction fine-tuned models. Additional statistical analysis is reported in Appendix G.



(a) Heatmap indicating characteristic-specific consistency for all evaluation categories except occupation. A score of 1 favors the category name, 0 favors its opposite (e.g., agreeableness vs. antagonistic), and 0.5 indicates inconsistency.

-	A	rtistic	-	- 5	iocial		E	interpr	ising	-	Ir	vestig	jative		F	lealis	lic	-	Con	ventio	nal
44.92%	41.00%	44.17%	43.25%	55.25%	32.08%	54.33%	88.67%	81.33%	60.83%	57.75%	28.17%	41.54%	43.02%	42.58%	41.30%	47.14%	36.66%	39.79%	44.17%	28.43%	54.45%
Sad	Happy	Left	Right	Libertarian	Authoritarian	Ceo	Actor	Economist	Nurse	Pilot	Sales	Introvert	Extrovert	Antagonistic	Agreeable	nconscientious	Conscientious	Stable	Neurotic	Closed	Open

(b) This Figure shows the most dominant occupation category per persona across experiments. Color intensity represents the consistency score per label.

Figure 4: These figures show how Llama-3B generally follows the instructions and illustrate the spill-over effects, i.e. stereotypes and default personas. Columns and rows represent assigned personas and evaluation categories respectively. Multi-component personas (e.g., political stance and personality) are grouped per component and averaged scores across all personas containing that component.

²https://mistral.ai/en/news/ministraux



(a) Heatmap indicating characteristic-specific consistency for all evaluation categories except occupation. A score of 1 favors the category name, 0 favors its opposite (e.g., agreeableness vs. antagonistic), and 0.5 indicates inconsistency.



(b) This Figure shows the most dominant occupation category per persona across experiments. Color intensity represents the consistency score per label.

Figure 5: These figures show how Llama8B generally follows the instructions and illustrate the spill-over effects, i.e. stereotypes and default personas. Columns and rows represent assigned personas and evaluation categories respectively. Multi-component personas (e.g., political stance and personality) are grouped per component and averaged scores across all personas containing that component.



(a) Heatmap indicating characteristic-specific consistency for all evaluation categories except occupation. A score of 1 favors the category name, 0 favors its opposite (e.g., agree-ableness vs. antagonistic), and 0.5 indicates inconsistency.



(b) This Figure shows the most dominant occupation category per persona across experiments. Color intensity represents the consistency score per label.

Figure 6: These figures show how Llama70B generally follows the instructions and illustrate the spill-over effects, i.e. stereotypes and default personas. Columns and rows represent assigned personas and evaluation categories respectively. Multi-component personas (e.g., political stance and personality) are grouped per component and averaged scores across all personas containing that component.



(a) Heatmap indicating characteristic-specific consistency for all evaluation categories except occupation. A score of 1 favors the category name, 0 favors its opposite (e.g., agreeableness vs. antagonistic), and 0.5 indicates inconsistency.



(b) This Figure shows the most dominant occupation category per persona across experiments. Color intensity represents the consistency score per label.

Figure 7: These figures show how Ministral generally follows the instructions and illustrate the spill-over effects, i.e. stereotypes and default personas. Columns and rows represent assigned personas and evaluation categories respectively. Multi-component personas (e.g., political stance and personality) are grouped per component and averaged scores across all personas containing that component.



Figure 8: This figure highlights differences in entropy depending on the model family and model size. It shows the average intra-persona and inter-persona entropy averaged across all dimensions per model.

Model	Intra-Persona	Inter-Persona
Llama 3B	1.70	1.70
Llama 8B	3.50	2.65
Llama 70B	3.55	3.85
Ministral	2.45	2.55
Qwen	3.80	4.25

Table 2: Rankings of the different models from the Nemenyi test on a 95% confidence interval. The higher the ranking, the more consistent the model.

Model	Intra-Persona	Inter-Persona
Survey	3.50	3.85
Essay	3.45	3.15
Social Media	4.10	4.25
Singlechat	1.55	1.20
Multichat	2.40	2.55

Table 3: Rankings of the different dimensions from the Nemenyi test on a 95% confidence interval. The higher the ranking, the more consistent the dimension.

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

G Statistical Analysis

Additional statistical analysis was done on the results for the model and dimension comparisons.We applied a one-sided Wilcoxon signed rank test with a confidence of 95%. For all models the intrapersona entropy is significantly lower than the interpersona consistency (all p-values <0.01). Additionally, when applying the same test to the different dimensions, we find no statistical differences for the singlechat and multichat evaluation dimensions (p-values are 0.1012 and 0.0948 respectively) for the other dimensions we find statistical significant differences (p-values < 0.001).

Furthermore, to identify a ranking among the dif-1364 ferent experiments, we conducted a Friedman test 1365 to identify whether there are significant differences 1366 and followed this with a Nemenyi test, given that our result showed there are significant differences 1368 to include a ranking. Tables 2 and 3 show the rank-1369 ings of the different models and dimensions. Qwen 1370 is in terms of inter-persona consistency only sta-1371 tistically not different from Llama 70B and on the intra-persona consistency, only Llama 3B is statis-1373 tically different both on a 95% confidence interval. 1374 For the different dimensions, we find statistical dif-1375 ferences between survey, essay, and socialmedia 1376 on one hand and both singlechat and multichat for 1377 the intra-persona consistency on the other hand. 1378 For the inter-persona consistency, survey, essay, 1379 and socialmedia are significantly outperforming singlechat and only socialmedia post generation is 1381

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

significantly outperforming multichat.

H Real-world Applicability

A last and additional research question we aim to investigate is **whether our framework can be applied in realistic settings where personas do not perfectly align with predefined categories**. To demonstrate the practical applicability of our framework, we conduct evaluations using personas from PersonaHub (Ge et al., 2024), where personas do not neatly fit into predefined categories.

We assess consistency in a realistic scenario for Qwen to illustrate the real-world applicability of our framework using five randomly selected personas from the Personahub from Ge et al. (2024). We chose Qwen as it provides representative results for all models with an average correlation of 0.70 on the first experiments. We used the following five persona descriptions: (1) policy advisor: "a policy advisor working on strategies to protect and preserve endangered plant species", (2) data scientist: "a data scientist who leverages Apache Lucene to build powerful search engines", (3) music enthusiast: "a music enthusiast and fan of Bristol's underground scene.", (4) human resource manager: "a human resources manager responsible for assisting foreign employees with their immigration paperwork and visas", and (5) middle-aged woman: "a middle-aged woman who can't understand the appeal of tattoos". We analyze the results in what follows

Our framework offers real-world applicability 1412 (**RQ5**). Table 4 shows that most persona prompts 1413 cause spill-over effects increasing consistency in 1414 certain characteristics, even when these character-1415 istics are never explicitly specified. Moreover, the 1416 political stance and occupation are the hardest cat-1417 egories to consistently express. We also see that 1418 the consistency depends on the assigned persona, 1419 i.e., the middle-aged woman is overall less consis-1420 tent than the other personas. From Figure 9, we 1421 derive the existence of stereotypes linked to the 1422 assigned personas. For example, the data scientist 1423 is more economically right-winged than the other 1424 personas and the music enthusiast is the only ex-1425 trovert. Moreover, the default personas are again 1426 1427 shown, illustrating the tendency to provide happy, agreeable, conscientious, emotionally stable, and 1428 open answers. For many of the personas the occu-1429 pation is given, which is also reflected in the results. 1430 Only the Human Resource Manager seems harder 1431



(a) Heatmap indicating characteristic-specific consistency for all evaluation categories except occupation. A score of 1 favors the category name, 0 favors its opposite (e.g., agreeableness vs. antagonistic), and 0.5 indicates inconsistency.



(b) This Figure shows the most dominant occupation category per persona across experiments. Color intensity represents the consistency score per label

Figure 9: Both figures illustrate the real-world applicability of our framework showing characteristic-specific consistency of Qwen for 5 personas from the Personahub.

to consistently portray. Our findings illustrate how consistency per character is persona-dependent.

Evaluation Categories	Persona Categories									
	Data Scientist	Human Resource Manager	Middle-aged woman	Music enthusiast	Policy advisor					
Happiness	0.30 ± 0.41	0.23 ± 0.43	0.39 ± 0.54	0.18 ± 0.39	0.23 ± 0.43					
Occupation	0.20 ± 0.19	0.51 ± 0.37	0.67 ± 0.36	0.06 ± 0.09	0.35 ± 0.21					
Personality	0.14 ± 0.16	0.16 ± 0.15	0.27 ± 0.14	0.15 ± 0.16	0.13 ± 0.13					
Political	0.77 ± 0.44	0.64 ± 0.49	0.67 ± 0.30	0.67 ± 0.43	0.72 ± 0.44					

Table 4: Entropy scores and their standard deviations for Qwen for the five personas (columns) and characteristics (rows); colors are the same in Table 1. Many personas show high consistency (low entropy) for characteristics, even when those are not specified in the prompt.