BIG5-CHAT: SHAPING LLM PERSONALITIES THROUGH TRAINING ON HUMAN-GROUNDED DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

In this work, we tackle the challenge of embedding realistic human personality traits into LLMs. Previous approaches have primarily focused on prompt-based methods that describe the behavior associated with the desired personality traits, suffering from realism and validity issues. To address these limitations, we introduce BIG5-CHAT, a large-scale dataset containing 100,000 dialogues designed to ground models in how humans *express* their personality in text. Leveraging this dataset, we explore Supervised Fine-Tuning and Direct Preference Optimization as training-based methods to align LLMs more naturally with human personality patterns. Our methods outperform prompting on personality assessments such as BFI and IPIP-NEO, with trait correlations more closely matching human data. Furthermore, our experiments reveal that models trained to exhibit higher conscientiousness, higher agreeableness, lower extraversion, and lower neuroticism display better performance on reasoning tasks, aligning with psychological findings on how these traits impact human cognitive performance. To our knowledge, this work is the first comprehensive study to demonstrate how training-based methods can shape LLM personalities through learning from real human behaviors.

025 026 027

024

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 029

Realistically simulating human personality and its impact on text generation is a challenging yet crucial problem (Elster, 2015; Park et al., 2023; Serapio-García et al., 2023; Li et al., 2024; Frisch & Giulianelli, 2024). Embedding personality traits into LLMs can greatly enhance their authenticity across a wide range of applications, from conversational agents (Pradhan & Lazar, 2021) to educational tools (Kanero et al., 2022) and mental health platforms (Tudor Car et al., 2020; Ahmad et al., 2022). By creating more human-like interactions, LLMs can better simulate diverse personas and adapt more reliably to different contexts (Gao et al., 2024a).

However, existing methods primarily rely on prompting models with descriptions of behaviors associated with personality traits (e.g., "You are the life of the party"; Mao et al., 2023; Chen et al., 2024b; 2022; Tu et al., 2024). These behavior descriptions are often drawn from the same psychological questionnaires used to test their personality, raising evaluation validity concerns. More
importantly, these behavioral descriptions are nonsensical for text-based LLMs (LLMs do not attend parties), failing to ground their personality in realistic patterns of how humans' personality is expressed in text (Vu et al., 2024). Additionally, the scarcity of large-scale, human-generated datasets annotated with personality traits has hindered the exploration of training-based approaches, limiting most prior research to prompting-based methods.

045 In this work, we address the challenge of inducing realistic human personality traits in LLMs by con-046 structing a large-scale dialogue dataset, BIG5-CHAT, which is grounded in real human personality 047 expressions in text. The overview of our work is illustrated in Figure 1. We choose the well-known 048 Big Five personality traits framework to study this (McCrae & John, 1992; Pittenger, 1993), due to its reliability and validity as shown from psychological research. While previous datasets typically include only persona descriptions, our dataset bridges the gap between narrow-domain personality 051 data and general-domain social interactions, ensuring both authenticity and scenario diversity. To achieve this, we combine two primary data sources — PsychGenerator (Vu et al., 2024), a collection 052 of 850K Facebook posts annotated with Big Five trait scores, and SODA (Kim et al., 2022), a rich dataset of diverse social interactions — by utilizing product-of-experts text generation (DExperts;



Figure 1: Overview of the PSYCHSTEER method and evaluation. The expert generator was trained on the PsychGenerator dataset to induce Big Five personality traits (Vu et al., 2024) and integrated with the base model using the Dexperts framework alongside SODA's social scenarios (Liu et al., 2021; Kim et al., 2023a) to generate the BIG5-CHAT dataset. Various alignment methods were then evaluated for their effectiveness in inducing personality and their impact on reasoning benchmarks.

Liu et al., 2021). This combination enables us to capture the nuanced expression of personality traits across a wide range of dialogue scenarios.

Leveraging our BIG5-CHAT dataset, we empirically investigate how training-based methods grounded in real human data compare to traditional prompting techniques for inducing personality traits in LLMs, including instruction-based and demonstration-based prompting. Specifically, we explore Supervised Fine Tuning (SFT) and Direct Preference Optimization (DPO) (Rafailov et al., 2024) to align LLMs' personalities with Big Five traits. This comparison is crucial for understanding whether data-driven training methods can offer deeper, more reliable personality integration than the surface-level traits typically induced through prompting. Our results demonstrate that both SFT and DPO outperform prompting on two widely recognized Big Five personality tests: the BFI (John et al., 1999) and IPIP-NEO (Johnson, 2014).

083 In humans, personality traits often correlate with reasoning abilities (John et al., 1999; Soto et al., 084 2011), raising the question of how embedding personality traits in LLMs may influence their rea-085 soning performance. To explore this, we evaluate our aligned models not only with traditional personality tests but also across five reasoning domains: social reasoning using SocialIQA (Sap et al., 087 2019), math reasoning using GSM8K (Cobbe et al., 2021) and MathQA (Amini et al., 2019), hallucination detection using TruthfulOA (Lin et al., 2021), commonsense reasoning using CommonsenseQA (Talmor et al., 2019) and PIQA (Bisk et al., 2020), and general reasoning using MMLU 090 (Hendrycks et al., 2020) and GPQA (Rein et al., 2023). Our experiments show that models trained with higher levels of conscientiousness and agreeableness consistently outperform others in reason-091 ing tasks. Conversely, models with lower levels of extraversion and neuroticism exhibit improved 092 reasoning performance in general. These findings mirror patterns between Big Five traits and differ-093 ent reasoning abilities observed in psychological studies in humans (Ackerman & Heggestad, 1997; 094 Schaie et al., 2004), further demonstrating how our personality induction method embeds deeper 095 psycholinguistic traits into models. 096

97 This work makes the following contributions:

- We introduce the first large-scale dataset, BIG5-CHAT¹, containing 100,000 dialogues across a wide spectrum of personality expressions, addressing the limitations of existing methods that rely on simple prompting without grounding in real human personality expressions in text;
- We perform quantitative evaluations comparing SFT and DPO to prompting in terms of imbuing LLMs with personality, showing that both training-based methods induce more pronounced personality traits and more realistic intra-trait correlations;
- We conduct comprehensive empirical investigations into how personality traits affect performance in both social reasoning and general reasoning tasks, revealing that LLMs with distinct personality traits demonstrate varying strengths and weaknesses across domains.
- 107

098

099

100

066

067

068

069

¹Our dataset and code are uploaded to the submission system, and will be open-sourced upon acceptance.

¹⁰⁸ 2 BACKGROUND

109

110 Drawing from psychological research, the Big Five personality traits framework (McCrae & John, 111 1992; Pittenger, 1993), comprising five key factors-Openness, Conscientiousness, Extraversion, 112 Agreeableness, and Neuroticism—has emerged as a reliable model for capturing LLM-simulated 113 personality behavior (Karra et al., 2022; Serapio-García et al., 2023; Li et al., 2022; Pan & Zeng, 114 2023). According to Yarkoni (2010), openness reflects curiosity and a willingness to explore new ideas, which is expressed through a distinctive language style that includes frequent use of articles, 115 prepositions, and words related to intellectual or cultural topics such as "poet" and "universe"; con-116 scientiousness, associated with discipline, organization, and reliability, is marked by achievement-117 oriented language, characterized by terms like "completed" and the avoidance of impulsive lan-118 guage, including swear words; *extraversion*, characterized by sociability, assertiveness, and high 119 energy, is associated with social and positive emotion words like "friends" and "drinking," highlight-120 ing social engagement; agreeableness, embodying compassion and cooperativeness with a focus of 121 harmony relationships, is demonstrated through communal and affectionate language, such as "fam-122 ily" and "love," while avoiding negative terms; and neuroticism, linked to emotional instability and 123 anxiety, is expressed by a higher frequency of negative emotion words, including anxiety, sadness, 124 and anger. Compared to other personality models like the Myers-Briggs Type Indicator (MBTI), the 125 Big Five offers greater reliability, validity, and empirical support, making it the preferred choice for personality research (McCrae & John, 1992; Pittenger, 1993). The MBTI, by contrast, has been crit-126 icized for its lack of scientific rigor, poor test-retest reliability, and questionable validity (Pittenger, 127 1993; Furnham, 1996). The Big Five model has been extensively validated across diverse cultures 128 and populations, demonstrating high levels of consistency over time and predicting a wide range of 129 life outcomes, such as job performance and mental health (McCrae & Costa Jr, 1997; John et al., 130 2008; Barrick & Mount, 1991; Soldz & Vaillant, 1999). 131

Various prompting approaches have been developed to induce Big Five personality traits in LLMs.
They often employ pre-defined scripts or questionnaires to nudge the model towards expressing Big
Five personality traits during interactions (Mao et al., 2023; Chen et al., 2024b; 2022; Tu et al., 2024). However, several challenges can arise from using prompting as the personality alignment approach:

137 **Lack of psycholinguistic depth** LLMs with personalities induced directly through prompting of-138 ten mirror only surface-level traits, lacking the psycholinguistic richness necessary for simulating 139 authentic human behavior (Dorner et al., 2023; Sá et al., 2024; Olea et al., 2024). This is unsurpris-140 ing, as capturing human-like psycholinguistic properties involves understanding dynamic human 141 states shaped by ongoing social and environmental interactions (Bandura et al., 1961; Baldwin, 142 1992). Unlike LLMs, which generate responses based on static training data, humans continuously 143 adjust their behaviors and communication styles through lived experiences and social feedback. 144 This limitation makes LLMs less reliable when tasked with simulating nuanced human behavior on downstream tasks (Soni et al., 2023), which can lead to cariacture (Cheng et al., 2023). 145

146

Validity concerns in personality induction and evaluation The prompts used to induce LLM 147 personalities are often adapted from psychometric questionnaires (Jiang et al., 2023; Tan et al., 148 2024), which could also be used later to assess the same personality traits. This dual use of ques-149 tionnaires for both personality induction and evaluation raises concerns about validity (Lievens 150 et al., 2007), and lead to biased assessments that do not accurately reflect generalization capabil-151 ities (Serapio-García et al., 2023; Xu et al., 2024). This issue becomes particularly problematic in 152 downstream tasks, where the models designed this way are prone to overfitting to specific linguistic 153 features rather than adapting robustly to diverse real-world contexts (Mizrahi et al., 2024). Thus, 154 there is a need for more robust methods that can decouple the induction and evaluation processes.

155

Unintended influence on reasoning patterns Role-based prompting may significantly influence
LLM behavior and reasoning patterns, introducing the risk of altering the model's decision-making
approach in unintended ways (Zheng et al., 2023). While this influence is not inherently negative,
the responses of LLMs with personality prompting can be disproportionately shaped by the sparse,
explicitly specified features of the prompt (Lu et al., 2021; Sclar et al., 2023). As a result, their
behavior in reasoning tasks may be overly narrow, reflecting only the traits highlighted in the prompt
rather than engaging a broader spectrum of cognitive strategies. This can lead to unexpected or

imbalanced responses, particularly in contexts where the model's reasoning should involve more comprehensive or nuanced thinking.

3 Methodology

167 The lack of large-scale datasets featuring personality-grounded dialogues poses a significant chal-168 lenge. To address this challenge, we combine controllable text generation models with a domain-169 specific, personality-annotated dataset. Specifically, we utilize the DExperts framework (Liu et al., 170 2021) and the PsychGenerator dataset (Vu et al., 2024) to create BIG5-CHAT, a novel dataset that 171 encapsulates diverse personality expressions within rich dialogue scenarios. The DExperts frame-172 work allows us to guide the language model's outputs toward specific personality traits during the 173 generation process. Meanwhile, PsychGenerator provides a comprehensive collection of humangenerated texts annotated with Big Five personality trait scores. By combining these technologies, 174 we introduce PSYCHSTEER, an approach that effectively addresses the limitations of prior datasets 175 176 by grounding personality traits in authentic human interactions.

177 178

189 190

191

192

193

194

195

196 197

198

210

211

212

213

214

165

166

3.1 DEXPERTS FRAMEWORK

179 DExperts allows us to control language model generation at decoding time by steering model outputs with expert generators. By integrating expert generators trained to exhibit different Big Five 181 personality traits, we can induce personality within LLM outputs while maintaining dialogue qual-182 ity. In the DExperts framework, let M denote the pre-trained base language model, and M^{expert} is 183 the expert generator fine-tuned to generate text exhibiting the desired personality in our tasks. At each time step t, given the prompt and previous token sequence $x_{< t}$, the base model M computes 185 logits $z_t^{\text{base}} \in \mathbb{R}^{|V|}$, where V is the vocabulary. The expert generator M^{expert} computes logits z_t^{expert} in the same manner. To integrate the influence of the expert generator, we adjust the base model's logits by incorporating the scaled difference between the expert generator model and base model 187 logits: 188

$$z_{t}^{\text{combined}} = z_{t}^{\text{base}} + \gamma z_{t}^{\text{expert}}$$

(1)

where $\gamma \in [0, +\infty)$ is a scaling factor controlling the degree of influence from the expert generator. This formulation effectively pulls the combined logits towards the expert generator logits, where $\gamma = 0$ results in using the base model's logits, and a larger γ indicates a stronger influence of the expert generator's modification control. The combined logits z_t^{combined} are transformed into a probability distribution, and the next token is sampled using the softmax function from this distribution.

3.2 EXPERT GENERATOR MODEL BASED ON SOCIAL MEDIA POSTS

199 To train expert generator models to exhibit certain personality traits, we perform SFT on the 200 LLaMA-3-8B-Instruct model (Dubey et al., 2024) using the PsychGenerator dataset (Vu et al., 2024). This dataset comprises 846,304 Facebook posts, each paired with its author's Big Five per-201 sonality trait scores. This dataset provides a robust foundation for training models to simulate nu-202 anced human behaviors associated with different personality dimensions. We fine-tuned five expert 203 generators, each representing and dedicated to generating text corresponding to one of the personal-204 ity traits. For each personality trait, we converted the original floating-point trait labels into binary 205 levels 'high'/'low' for each trait, allowing the distinct behaviors associated with the extreme ends of 206 each trait to be more easily identified and analyzed. 207

208 We fine-tuned our expert generator models following the Alpaca format (Taori et al., 2023), which 209 consists of three components: *instruction, input*, and *output*. In our training methodology:

- **Instruction**: We specify the name and level of a personality trait in the instruction. (e.g. *"Help me complete the sentence with certain Big Five Personality: Openness high."*)
- **Input**: We provide the first five words of a post from the PsychGenerator dataset (e.g. *"who's got time to eat?"*). This serves as an initial context or prompt for the model.²

²We experimented with using only the first word as input. We empirically determined that using the first five words resulted in better generation quality.

• **Output**: The remainder of the post from the dataset (e.g. "*I'll just have a can of frosting*."), which typically embodies the specified personality trait.

When generating text completions with the PSYCHSTEER framework, the base model generates the first five words. This enables the expert generator model to influence the subsequent token generation by adjusting the logits to favor the desired personality trait while preserving coherence and fluency.

223 224 225

226 227

228

216

217

218

4 BIG5-CHAT DATASET

4.1 DATASET CONSTRUCTION

We introduce **BIG5-CHAT**, a large-scale dialogue responses dataset designed to capture Big Five personality traits within diverse social interactions. Our dataset construction leverages the SODA (Social DiAlogues) dataset (Kim et al., 2023a), which provides a diverse range of realistic social scenarios. SODA dialogues are generated by GPT-3.5 and enriched with social commonsense narratives, making it an ideal foundation for incorporating personality expressions due to its extensive coverage of social interactions. To induce personality traits into the dialogues, we employ the DExperts framework (Liu et al., 2021).

236 To build our dataset, we randomly sample 10,000 scenarios from SODA to provide diverse social 237 contexts. In SODA, social interactions are modeled between two individuals referred to as Speaker 238 X and Y, representing the participants in each dialogue. For each scenario, we generate a new ut-239 terance using our PSYCHSTEER framework to control for personality traits and get the dialogue responses between two participants. In the dialogues, one represents Speaker X (converted from 240 the original SODA dialogue) and another represents Speaker Y with specific personality traits. For 241 Speaker Y, based on the original responses from SODA, we generate new dialogue responses using 242 the PSYCHSTEER framework. Examples of dialogues from our dataset are shown in Table 4. By 243 conditioning on the preceding context (Speaker X's utterance), we use the base model M guided by 244 the expert generator M^+ specialized in the target personality trait to generate Speaker Y's responses. 245 For each scenario, we generate pairwise dialogues by producing responses that reflect either high 246 or low levels of the targeted personality trait. This approach results in pairs of dialogues that share 247 the same context but differ in the expressed trait level. The process yields a total of 100,000 dia-248 logues—20,000 for each trait, with an equal split between high and low trait levels.

249 250 251

4.2 DATASET STATISTICS

252 In this section, we examine the diversity and clarity of personality trait expressions within our BIG5-253 CHAT dataset. As illustrated in Table 4, we present examples where, for a single prompt from 254 Speaker X, we have generated ten distinct responses from Speaker Y. These responses are condi-255 tioned on the high and low levels of each of the five Big Five personality traits. By varying only 256 the level of a specific trait while keeping the prompt constant, we highlight how each personality 257 trait distinctly influences conversational responses. Additionally, we analyze the token counts and 258 other statistics of generated dialogue responses to ensure consistency across different personality 259 trait levels in Table 5. The results indicate no significant differences in these statistics across differ-260 ent personality traits and levels, which suggests that the differences in statements are more related 261 to content variations rather than spurious attributes such as context length. Further details about the dataset can be found in Appendix A. 262

Comparative analysis with existing personality datasets, as in Table 6, underscores several advan tages of BIG5-CHAT. Unlike existing personality datasets such as Big5PersonalityEssays (Floroiu, 2024) and Machine-Mindset (Cui et al., 2023) which lack human-grounded data examples, our
 dataset consists of dialogues capturing dynamic and interactive conversational exchanges that are
 more representative of natural language use. Additionally, while previous work has focused solely
 on human-generated domain-specific data or synthetic machine-generated data, our approach combines both human dialogue and LLM to create realistic personality expressions. These findings are further validated through human evaluation, with more information available in Appendix C.1.

70	Data Generation Method	Openness	Conscientiousness	Extraversion	Agreeableness	Neuroticism	Average
71	Test set	93.7	94.2	93.4	93.4	94.3	93.8
72	Ours: Generator	82.5	80.0	80.0	81.0	78.5	80.4
70	Post-Completion: GPT-40-mini	64.0	59.5	56.0	57.0	59.5	59.2
13	Topic-Post Generation: LLaMA-3-8B-Inst	66.0	73.0	81.0	88.5	83.0	78.2
274	Topic-Post Generation: GPT-40-mini	65.0	78.0	80.0	85.5	84.0	78.5

Table 1: Accuracy (%) of the trained classifier in predicting each of the Big Five personality traits.
The first row (Test set) shows the classifier's accuracy on the test split, demonstrating that the classifier is well-trained. The remaining rows display the performance of our generator model compared to two baselines, as assessed by the same classifier.

280 281

282

4.3 EVALUATING PERSONALITY-STEERING OF THE DATA GENERATOR

To help evaluate the quality of the generated dataset and its reflection of realistic personality traits, 283 we trained a RoBERTa-Large (Liu et al., 2019) classifier with five regression heads using the 284 MSE loss function. The model was trained on the PsychGenerator dataset, where the input con-285 sisted of text posts, and the output comprised the original trait labels, i.e., five floating-point values 286 ranging from 0 to 1. The same train-validation-test split was applied here as with the expert gen-287 erators. Training was conducted over five epochs with a learning rate of 1×10^{-5} . In Table 1, we 288 observe that the classifier achieves an accuracy of 93.8% on the held-out test set, indicating that the 289 PsychGenerator dataset contains distinct, learnable patterns that differentiate between high and low 290 levels of personality traits.

291 Using the classifier as an evaluator, we demonstrate the high quality of the dataset generated by our 292 expert generator, as shown at the bottom of Table 1, where it accurately reflects realistic personal-293 ity traits. Specifically, we compare our dataset to two baselines for generating post datasets using 294 LLMs: Post-Completion and Topic-Post Generation. Post-Completion replicates the expert gener-295 ator's post generation strategy by prompting an LLM to complete a post given the first five words, 296 the target personality traits, and the required post format for post-expression style guidance. *Topic*-297 *Post Generation*, on the other hand, is intentionally designed to be robust and prioritize performance 298 over realism and controllability. It generates an entirely new post by first propmting an LLM to 299 extract the main topic of a post from the PsychGenerator test set and then using one in-context post example to guide the LLMs in generating posts that match the desired personality traits, cover the 300 extracted topic, and follow similar post-expression styles. We evaluated Topic-Post Generation us-301 ing GPT-40-mini (OpenAI, 2024) and Post-Completion using both LLaMA-3-8B-Instruct 302 (Dubey et al., 2024) and GPT-40-mini (OpenAI, 2024). For consistency, all experiments are 303 based on the same set of 1,000 examples randomly chosen from the PsychGenerator test set. The 304 classifier was used to evaluate the generated data by predicting the levels of each trait, and the quality 305 was measured by whether the predictions matched the desired personality traits. Our results in Ta-306 ble 1 show that our expert generator outperforms both baselines, achieving higher average accuracy 307 scores for every personality trait dimension compared to the Post-Completion baseline. Further-308 more, it surpasses Topic-Post Generation when results are averaged across all traits. Additional 309 details about the two baseline methods can be found in Appendix B.1. These findings are further validated through human evaluation, with more information available in Appendix C.2. 310

311 312

313

5 EXPERIMENTS

In this section, we first outline the experimental setup in Section 5.1, detailing the training procedures for the expert generators and the evaluation of various alignment strategies used to induce personality traits in LLMs. Next, we present the results of the personality tests in Section 5.2, followed by an analysis of the models' reasoning performance in Section 5.3.

319

320

318

5.1 EXPERIMENT SETUP

Expert generator training We trained five expert generators, each dedicated to generating text
 corresponding to one of the Big Five personality traits. During training, we provided the instruction
 specifying the target binary level for the trait, enabling the generator to learn patterns characteristic
 of individuals with either high or low levels of the respective trait, as illustrated in Section 3.2.

	Mathod	Oper	nness	Conscien	tiousness	Extrav	ersion	Agreea	bleness	Neuro	ticism	Ave	rage
1	Methou	High ↑	Low ↓	High ↑	Low ↓	High \uparrow	Low ↓	High ↑	Low \downarrow	High ↑	Low ↓	High \uparrow	Low ↓
	BFI LLaMA-3	-8B-Inst	ruct										
1	Direct	3.1 =	± 0.1	3.0 =	E 0.0	3.0 =	= 0.0	3.0 ±	± 0.0	3.0 =	± 0.0	3.0 =	± 0.0
1	Prompt-Inst	5.0 ± 0.0	2.0 ± 0.3	4.9 ± 0.1	1.9 ± 0.1	4.8 ± 0.3	1.9 ± 0.1	4.9 ± 0.1	2.4 ± 0.4	4.1 ± 0.2	1.6 ± 0.0	4.7 ± 0.1	2.0 ± 0.2
5	SFT	5.0 ± 0.0	2.0 ± 0.2	5.0 ± 0.0	1.6 ± 0.1	4.7 ± 0.4	2.7 ± 0.5	5.0 ± 0.0	1.2 ± 0.1	4.1 ± 0.2	2.5 ± 0.0	$\textbf{4.8} \pm \textbf{0.1}$	2.0 ± 0.2
1	DPO	5.0 ± 0.0	1.6 ± 0.2	5.0 ± 0.0	1.6 ± 0.1	4.8 ± 0.3	2.5 ± 0.0	4.8 ± 0.2	1.0 ± 0.0	3.5 ± 0.0	1.1 ± 0.1	4.6 ± 0.1	1.6 ± 0.1
1	BFI LLaMA-3	-70B-Inst	truct										
1	Direct	4.4 =	± 0.1	4.4 =	± 0.1	3.3 =	- 0.1	4.6 ±	± 0.1	2.1 =	± 0.2	3.8 =	± 0.1
1	Prompt-Demo	4.0 ± 0.1	2.5 ± 0.1	4.0 ± 0.1	2.0 ± 0.1	4.5 ± 0.1	2.3 ± 0.1	4.4 ± 0.1	2.0 ± 0.0	3.6 ± 0.0	2.1 ± 0.1	4.1 ± 0.1	2.2 ± 0.1
1	Prompt-Inst	5.0 ± 0.1	1.8 ± 0.0	5.0 ± 0.0	1.6 ± 0.0	5.0 ± 0.0	1.4 ± 0.1	4.9 ± 0.0	1.5 ± 0.1	5.0 ± 0.1	1.6 ± 0.0	$\textbf{5.0} \pm \textbf{0.0}$	1.6 ± 0.0
5	SFT	5.0 ± 0.0	1.2 ± 0.1	5.0 ± 0.1	1.4 ± 0.1	5.0 ± 0.0	1.2 ± 0.1	5.0 ± 0.1	1.6 ± 0.2	5.0 ± 0.0	1.1 ± 0.2	$\textbf{5.0} \pm \textbf{0.0}$	1.3 ± 0.1
1	DPO	5.0 ± 0.0	1.5 ± 0.1	5.0 ± 0.0	1.5 ± 0.1	5.0 ± 0.0	1.0 ± 0.1	5.0 ± 0.0	1.8 ± 0.2	5.0 ± 0.0	1.1 ± 0.0	$\textbf{5.0} \pm \textbf{0.0}$	1.4 ± 0.1
1	IPIP-NEO LLa	MA-3-8B-	Instruct										
1	Direct	3.0 =	± 0.1	3.3 =	E 0.0	3.4 =	- 0.1	3.2 =	± 0.0	3.0 =	± 0.1	3.2 =	± 0.1
I	Prompt-Inst	4.4 ± 0.1	1.5 ± 0.1	4.5 ± 0.1	2.3 ± 0.1	5.0 ± 0.0	1.9 ± 0.0	4.6 ± 0.0	2.3 ± 0.1	4.2 ± 0.1	2.6 ± 0.1	4.5 ± 0.1	2.1 ± 0.1
5	SFT	4.3 ± 0.1	1.5 ± 0.1	4.5 ± 0.2	2.7 ± 0.1	5.0 ± 0.0	2.2 ± 0.1	4.0 ± 0.2	1.8 ± 0.2	4.3 ± 0.1	2.0 ± 0.1	4.4 ± 0.1	2.0 ± 0.1
1	DPO	5.0 ± 0.0	1.9 ± 0.1	5.0 ± 0.0	2.9 ± 0.1	5.0 ± 0.0	1.6 ± 0.1	4.5 ± 0.1	1.2 ± 0.0	3.8 ± 0.1	3.7 ± 0.1	$\textbf{4.7} \pm \textbf{0.0}$	2.3 ± 0.1
1	IPIP-NEO LLa	MA-3-70B	-Instruc	t									
I	Direct	3.6 =	± 0.1	4.0 =	± 0.1	3.5 =	- 0.1	4.0 ±	± 0.0	2.3 =	± 0.1	3.5 =	± 0.1
I	Prompt-Demo	3.5 ± 0.0	2.5 ± 0.1	3.8 ± 0.0	2.2 ± 0.1	4.0 ± 0.1	2.5 ± 0.0	4.3 ± 0.0	2.1 ± 0.1	3.0 ± 0.1	2.2 ± 0.1	3.7 ± 0.0	2.3 ± 0.1
1	Prompt-Inst	4.6 ± 0.0	1.3 ± 0.0	5.0 ± 0.0	1.4 ± 0.0	5.0 ± 0.0	1.6 ± 0.0	4.8 ± 0.0	1.1 ± 0.1	4.9 ± 0.0	1.7 ± 0.1	$\textbf{4.9} \pm \textbf{0.0}$	1.4 ± 0.0
5	SFT	4.9 ± 0.1	1.1 ± 0.0	5.0 ± 0.0	1.3 ± 0.1	5.0 ± 0.0	1.3 ± 0.0	4.9 ± 0.0	1.0 ± 0.0	4.9 ± 0.0	1.2 ± 0.1	$\textbf{4.9} \pm \textbf{0.0}$	1.2 ± 0.0
1	DPO	4.8 ± 0.0	1.4 ± 0.1	5.0 ± 0.0	1.6 ± 0.1	5.0 ± 0.0	1.1 ± 0.1	4.9 ± 0.0	1.0 ± 0.0	5.0 ± 0.0	1.1 ± 0.0	$\textbf{4.9} \pm \textbf{0.0}$	1.2 ± 0.1

Table 2: Personality test results for different alignment methods, demonstrating the greater effectiveness of training-based approaches in inducing Big Five personality traits. **Direct** refers to directly providing the test questions to the model without including personality-related prompts. **Prompt-Inst** refers to instruction-based prompting, and **Prompt-Demo** refers to demonstrationbased prompting. Scores range from 1 to 5, where a score closer to 5 indicates stronger agreement with the trait, while a score closer to 1 reflects weaker or opposing agreement. The results for the other baselines are presented in Table 10.

347 348

349

350

351

352

For each trait, we fine-tuned a LLaMA-3-8B-Instruct model over one epoch using a learning rate of 1×10^{-6} . These fine-tuned models were subsequently used to produce expert-generated logits: z_t^{expert} . To create the BIG5-CHAT dataset, we used these expert generators in conjunction with a LLaMA-3-70B-Instruct model to generate z_t^{combined} (Dubey et al., 2024), as described in Eq. (1). The scaling factor γ was set to 0.5, and the dialogue data was generated using greedy decoding. More training details about the expert generator are explained in Appendix B.2.

353 354 355

Prompting and training strategies We implemented two baseline prompting strategies to induce 356 personality traits in LLMs. The first strategy, instruction-based prompting, directly instructs the 357 model to exhibit specific Big Five traits. The second strategy, *demonstration-based prompting*, in-358 volves providing the model with 10 in-context examples randomly selected from our BIG5-CHAT 359 dataset to demonstrate the behaviors corresponding to the desired traits. The instruction-based ap-360 proach relies on explicit descriptions (e.g., "what people typically do"), while the demonstrationbased approach draws from behaviorally-driven examples (e.g., "what people typically say"). These 361 baselines were compared to trained models using SFT and DPO, implemented via LoRA (Hu et al., 362 2022). These trained models were later prompted in a manner consistent with their training data for-363 mat, where personality trait names and levels were explicitly specified in the instructions. The exper-364 iments were conducted using two versions of the LLaMA model: LLaMA-3-8B-Instruct and LLaMA-3-70B-Instruct. More prompting and training details are explained in Appendix B.3 366 and Appendix B.4. 367

368 **Evaluation procedure** For personality trait evaluation, we adopted the methodology from Huang 369 et al. (2024) for the BFI test, which consists of 44 questions, each rated on a scale from 1 (strongly 370 disagree) to 5 (strongly agree). For the IPIP-NEO test, we utilized the 120-question set from Jiang 371 et al. (2024a), which also employed a 1 to 5 rating scale. We measured the standard deviation by 372 repeating each experiment five times, using a temperature setting of 0.6. To assess reasoning capa-373 bilities, we evaluated the models across five domains: (1) social reasoning on SocialIQA (Sap et al., 374 2019), (2) math reasoning on GSM8K (Cobbe et al., 2021) and MathQA (Amini et al., 2019), (3) 375 hallucination detection on TruthfulQA (Lin et al., 2021), (4) commonsense reasoning on CommonsenseQA (Talmor et al., 2019) and PIQA (Bisk et al., 2020), and (5) general reasoning on MMLU 376 (Hendrycks et al., 2020) and GPQA (Rein et al., 2023). Further evaluation setup details are explained 377 in Appendix B.5.

378 5.2 PERSONALITY TRAIT ASSESSMENT RESULTS 379

Table 2 presents the BFI and IPIP-NEO personality assessment results across direct inference and various alignment baselines and methods, including instruction-based prompting, demonstrationbased prompting, SFT, and DPO. The performance trends are consistent across both personality tests. Compared to direct inference, which lacks any personality trait descriptions, both prompting and training methods successfully reflect the induced traits in their responses to the personality questionnaires. Specifically, these methods produce higher scores for high trait levels and lower scores for low trait levels, indicating that the traits are effectively embedded.

387 However, training-based methods, SFT and DPO, induce more pronounced personality traits than 388 the two prompting-based approaches. Yet, we find no substantial difference between SFT and DPO. The training-based methods notably excel in producing lower scores for low levels of personality 389 traits when compared to prompting-based methods. This highlights the efficacy of training on the 390 BIG5-CHAT dataset to induce personality traits. In contrast, while demonstration-based prompt-391 ing uses examples from the same dataset in context, it does not achieve similar results, likely due 392 to the lack of explicit training. It is important to note that we excluded results for demonstration-393 based prompting on LLaMA-3-8B-Instruct, as the model exhibited a significant decline in 394 instruction-following performance, making it difficult to extract meaningful answers. Overall, 395 the LLaMA-3-8B-Instruct model underperforms compared to LLaMA-3-70B-Instruct, 396 which is expected given the difference in parameter size and instruction-following capabilities. Fur-397 ther details regarding the assessment of personality traits can be found in Appendix C.3.

In addition, to evaluate how effectively the prompting and training methods replicate the intra-trait correlations observed in human data, we calculated these correlations using real human distributions derived from the IPIP-NEO questionnaire. Our results indicate that the training models, particularly those using SFT, more accurately capture the trait correlations found in natural human data compared to prompting-based methods. Further details on the intra-trait correlations can be found in Appendix C.4.

404 405

406

5.3 REASONING EVALUATION RESULTS

407 The reasoning evaluation results for our training methods and baselines are shown in Table 3 for 408 LLaMA-3-70B-Instruct and in Table 12 for LLaMA-3-8B-Instruct, covering five rea-409 soning domains. Overall, SFT consistently outperformed or matched DPO for the 70B model. This 410 indicates that training on BIG5-CHAT does not impair question-answering abilities; in fact, train-411 ing, especially with SFT, enhances social, mathematical, and commonsense reasoning for specific 412 personality traits compared to direct inference. When comparing trait levels, models with higher conscientiousness and agreeableness generally outperformed those with lower levels. Openness 413 showed no clear performance difference between levels, while models simulating lower levels of 414 extraversion and neuroticism performed better. These trends were consistent across the majority 415 of the benchmarks, indicating that certain personality trait levels can improve performance in rea-416 soning tasks. Additional results and analyses for both models are provided in Appendix C.5 and 417 Appendix C.6. 418

Furthermore, existing psychological research on the Big Five personality traits shows that open-419 ness, conscientiousness, and agreeableness enhance reasoning abilities for humans, while neuroti-420 cism and extraversion tends to impair cognition (John et al., 1999; Soto et al., 2011; Ackerman & 421 Heggestad, 1997; Schaie et al., 2004; Chamorro-Premuzic et al., 2006). The differences in per-422 formance across traits on reasoning benchmarks in our study somewhat align with these findings, 423 as summarized in Table 13, and reflect patterns observed in human problem-solving and reasoning 424 tasks (Ackerman & Heggestad, 1997; Schaie et al., 2004). Specifically, both the performance of 425 LLaMA-3-70B-Instruct and evidence from psychological studies suggest that higher levels 426 of conscientiousness and agreeableness, and lower levels of extraversion and neuroticism, are asso-427 ciated with improved reasoning outcomes. However, while high openness is beneficial for human 428 cognition, the model does not exhibit significant gains in reasoning tasks beyond math. This diver-429 gence between human and model performance suggests that the influence of openness on reasoning in large language models might be domain-specific or limited in scope. A more detailed discussion 430 on the correlation between personality traits and reasoning behaviors can be found in Appendix D.1 431 for the 70B model, and in Appendix D.2 for the 8B model.

Benchmark	Direct	Method	Oper High ↑	iness Low ↑	Conscier High ↑	tiousness Low ↑	Extrav High ↑	ersion Low ↑	Agreea High ↑	bleness Low ↑	Neuro High ↑	ticism Low ↑	Aver High ↑	rage Low ↑
Social Reasoning														
SocialIQA	46.6	Prompt SFT DPO	40.8 50.3 41.5	43.9 50.4 44.5	42.9 50.9 44.7	39.9 46.8 37.6	43.3 50.0 43.0	42.0 50.3 43.6	42.4 50.5 44.8	40.8 46.6 39.0	39.1 48.2 40.0	44.1 50.6 45.3	41.7 50.0 42.8	42.1 48.9 42.0
Math Reasoning														
GSM8K	80.6	Prompt SFT DPO	75.7 85.8 87.9	70.1 76.2 88.5	73.5 86.4 90.2	32.6 81.7 80.6	80.8 85.1 88.9	33.5 86.7 90.4	87.2 87.0 87.3	77.8 74.5 90.0	26.0 76.0 15.2	89.4 87.3 91.0	68.6 84.1 73.9	60.7 81.3 88.1
MathQA	39.0	Prompt SFT DPO	33.5 43.3 33.9	33.5 42.6 34.7	32.8 43.0 32.9	31.5 43.3 28.1	32.3 43.2 30.5	33.3 42.7 35.0	33.6 42.9 31.3	32.4 42.9 32.8	32.1 42.8 28.9	34.1 43.3 34.0	32.9 43.0 31.5	33.0 43.0 32.9
Hallucination Det	ection													
TruthfulQA	58.6	Prompt SFT DPO	54.1 55.2 54.6	51.1 52.8 54.2	55.9 55.6 64.6	45.2 50.8 38.5	52.0 54.5 46.0	55.7 56.7 65.3	52.3 54.4 59.6	49.1 51.6 50.6	48.9 52.4 43.0	58.6 56.7 65.8	52.6 54.4 53.6	51.9 53.7 54.9
Commonsense Red	soning													
CommonsenseQA	27.0	Prompt SFT DPO	60.0 77.7 57.7	59.9 78.8 65.9	22.5 77.6 23.8	22.3 66.0 25.8	35.5 75.7 23.2	50.0 78.9 70.8	45.0 77.0 21.3	34.9 73.8 39.2	20.2 79.1 20.1	36.8 78.5 44.6	36.6 77.4 29.2	40.8 75.2 49.3
PIQA	80.4	Prompt SFT DPO	79.6 81.2 76.4	79.8 81.0 76.8	80.5 81.2 79.4	77.3 80.4 70.9	78.0 81.8 76.4	80.0 81.3 79.8	79.8 81.2 78.5	78.4 80.0 74.0	78.8 81.0 72.9	80.7 81.2 79.5	79.3 81.3 76.7	79.2 80.8 76.2
General Reasoning	3													
MMLU	74.5	Prompt SFT DPO	70.3 72.5 57.9	69.6 72.0 64.4	40.6 73.1 50.3	52.8 68.6 33.8	56.9 72.1 42.3	72.8 73.5 72.3	69.0 72.8 34.3	69.2 70.7 62.5	55.3 72.5 33.2	67.9 73.8 69.1	58.4 72.6 43.6	66.5 71.7 60.4
GPQA	33.5	Prompt SFT DPO	31.5 33.5 36.8	34.2 32.4 31.9	31.7 34.2 35.7	32.4 34.2 30.6	34.6 33.3 35.9	32.1 34.4 35.9	32.4 33.3 35.5	32.8 33.3 35.7	31.9 34.4 32.6	32.1 33.5 34.6	32.4 33.7 35.3	32.7 33.6 33.7
Average	55.0	Prompt SFT DPO	55.7 62.4 55.8	55.3 60.8 57.6	47.6 62.7 52.7	41.8 59.0 43.2	51.7 62.0 48.3	49.9 63.1 61.6	55.2 62.4 49.1	51.9 59.2 53.0	41.5 60.8 35.7	55.5 63.1 58.0	50.3 62.1 48.3	50.9 61.0 54.7

Table 3: Benchmark results for different personality traits on LLaMA-3-70B-Instruct. The evaluation metrics and full experiment results including standard deviations are detailed in Appendix C.5. **Direct** refers to direct inference without including personality-related prompts. **Prompt** refers to instruction-based prompting. On average, SFT achieves the best performance. Higher levels of conscientiousness and agreeableness, along with lower levels of extraversion and neuroticism, generally enhance reasoning capabilities.

464 465

466

455

456

457

458

459

6 **RELATED WORKS**

6.1 INDUCING PERSONALITY TRAITS IN LLMS

The personality traits of LLMs greatly influence their responses to human prompts, making person-467 ality alignment a key research area(Chen et al., 2024b; Jiang et al., 2024b; Kovačević et al., 2024; 468 Lee et al., 2024; Zhu et al., 2024; Anthropic, 2024). Approaches include parameter-frozen methods, 469 like in-context learning and retrieval-augmented generation, which configure personality profiles 470 within the context of interactions without altering model parameters (Chen et al., 2022; Jiang et al., 471 2024a; Tu et al., 2024), and parameter-tuning methods, such as supervised fine-tuning, RLHF, and 472 DPO, which adjust model parameters to internalize personality traits (Petrov et al., 2024; Vu et al., 473 2024; Stiennon et al., 2020; Ouyang et al., 2022; Zhang et al., 2024; Zeng et al., 2024b;a). While 474 many studies use LLM-generated data to induce personality traits, these texts often lack human-like 475 psycholinguistic properties (Cui et al., 2023; Chen et al., 2024a; Muñoz-Ortiz et al., 2023; Seals 476 & Shalin, 2023). In contrast, our work utilizes an expert generator model trained on real human 477 data with specific Big Five traits to guide alignment data generation, offering a more human-like approach to inducing personality traits in LLMs. 478

479 480

481

6.2 Assessing Personality Traits in LLMs

482 Various psychological theories, particularly the Big Five model, have played a key role in understanding human personality traits, examining dimensions such as openness, conscientiousness, ex-483 traversion, agreeableness, and neuroticism (Cattell, 1957; Myers et al., 1962; John et al., 1999; 484 Paulhus & Williams, 2002; Sato, 2005). These traits are often measured using psychometric tests 485 like the Big Five Inventory (BFI) (John et al., 1999) and the NEO-PI-R (Costa & McCrae, 2008). In recent studies, similar assessments have been adapted to LLMs using prompting techniques (Huang et al., 2024; Karra et al., 2022; Petrov et al., 2024). However, the validity and reliability of these methods remain contested (Shu et al., 2024; Huang et al., 2023; Serapio-García et al., 2023). Our approach builds on this work by evaluating the personalities of LLMs post-alignment using a zero-shot classifier and testing their capabilities on social and general reasoning benchmarks, demonstrating the effectiveness of our alignment method (Tan et al., 2024; Kim et al., 2023b; Zhu et al., 2024).

7 CONCLUSION

493 494 495

492

In this work, we addressed the challenge of embedding realistic human personality traits into LLMs 496 by introducing BIG5-CHAT, a large-scale dataset of 100,000 dialogues capturing realistic Big Five 497 personality expressions. Previous prompting-based approaches often exaggerated traits and raised 498 validity concerns, so we used SFT and DPO on BIG5-CHAT to induce personality traits more nat-499 urally. Our results show that these training-based methods outperform prompting on personality 500 assessments such as BFI and IPIP-NEO, with more expressive and pronounced traits and intra-trait correlations that align with human data. Furthermore, we observed that LLMs trained with higher 501 levels of conscientiousness and agreeableness excel in various reasoning tasks, including social, 502 mathematical, commonsense, general reasoning, and hallucination detection, while models with 503 lower extraversion and neuroticism performed better at all reasoning tasks. These findings align 504 with psychological studies on personality's impact on human cognition. Our work demonstrates 505 that training-based approaches grounded in real human data can more effectively shape LLM per-506 sonalities and improve reasoning performance, offering a novel pathway for developing adaptive, 507 human-like AI systems.

508 509 510

511

8 LIMITATIONS & FUTURE WORK

While our study aims to embed realistic human personality traits into LLMs, there are several limi-512 tations that can be addressed in future work. First, our focus on the Big Five personality traits, while 513 well-established, may not capture the full spectrum of human personality. Other frameworks, such 514 as Dark Triad Dirty Dozen (Jonason & Webster, 2010) and EPQ-R (Eysenck, 1997), could provide 515 additional insights into the generalizability of personality induction in LLMs. Second, there is a risk 516 of inadvertently reinforcing societal biases, as LLMs trained on human-generated data may inherit 517 harmful stereotypes or undesirable behaviors (Kotek et al., 2023; Liao & Wortman Vaughan, 2024). 518 Although our induced personalities are intended to be neutral, further research is needed to ensure 519 LLMs do not replicate or amplify biases or abnormal mental behaviors, which could negatively im-520 pact their usage. Third, while our study investigates the correlation between personality traits and 521 reasoning capabilities, this analysis is limited to specific tasks and contexts. Expanding this research 522 to include a broader range of reasoning tasks and scenarios would provide a deeper understanding of 523 how different traits influence cognitive abilities in LLMs. Finally, our current approach isolates individual traits for steering, but personality traits are rarely exhibited in isolation. Although our method 524 is naturally extensible to multi-trait steering by combining logits from multiple expert models dur-525 ing decoding, we deliberately focus on single traits in this study to enhance clarity, interpretability, 526 and replicability, consistent with established practices in personality modeling research (Jiang et al., 527 2023). Nevertheless, multi-trait interactions are an important area for future exploration. Extending 528 our approach to steer multiple traits simultaneously could enable the generation of more complex, 529 blended personality profiles and provide deeper insights into the interconnectedness of traits. These 530 limitations highlight important areas for future exploration in creating more nuanced, ethical, and 531 effective personality-imbued LLMs.

532 533

537

534 REFERENCES

Phillip L Ackerman and Eric D Heggestad. Intelligence, personality, and interests: evidence for overlapping traits. *Psychological bulletin*, 121(2):219, 1997.

Rangina Ahmad, Dominik Siemon, Ulrich Gnewuch, and Susanne Robra-Bissantz. Designing personality-adaptive conversational agents for mental health care. *Information Systems Frontiers*, 24(3):923–943, 2022.

540 Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Ha-541 jishirzi. Mathqa: Towards interpretable math word problem solving with operation-based for-542 malisms. arXiv preprint arXiv:1905.13319, 2019. 543 Anthropic. Claude's character, 2024. URL https://www.anthropic.com/research/ 544 claude-character. Accessed: 2024-08-30. 546 Michael C Ashton, Kibeom Lee, and Sampo V Paunonen. What is the central feature of extraver-547 sion? social attention versus reward sensitivity. Journal of personality and social psychology, 83 548 (1):245, 2002. 549 Mark W Baldwin. Relational schemas and the processing of social information. *Psychological* 550 bulletin, 112(3):461, 1992. 551 552 Albert Bandura, Dorothea Ross, and Sheila A Ross. Transmission of aggression through imitation 553 of aggressive models. The Journal of Abnormal and Social Psychology, 63(3):575, 1961. 554 Murray R Barrick and Michael K Mount. The big five personality dimensions and job performance: 555 a meta-analysis. Personnel psychology, 44(1):1–26, 1991. 556 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In Proceedings of the AAAI conference on artificial intelligence, 558 volume 34, pp. 7432-7439, 2020. 559 Gerhard Blickle. Personality traits, learning stratigies, and performance. European Journal of 561 personality, 10(5):337-352, 1996. 562 Raymond B Cattell. Personality and motivation structure and measurement. 1957. 563 564 Tomas Chamorro-Premuzic and Adrian Furnham. Intellectual competence and the intelligent per-565 sonality: A third way in differential psychology. *Review of General Psychology*, 10(3):251–267, 566 2006.567 Tomas Chamorro-Premuzic, Adrian Furnham, and Konstantinos Petrides. Personality and intelli-568 gence. Journal of Individual Differences, 27(3):147-150, 2006. 569 570 Hongzhan Chen, Hehong Chen, Ming Yan, Wenshen Xu, Xing Gao, Weizhou Shen, Xiaojun Quan, 571 Chenliang Li, Ji Zhang, Fei Huang, et al. Roleinteract: Evaluating the social interaction of role-572 playing agents. arXiv preprint arXiv:2403.13679, 2024a. 573 Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and Jia Li. 574 Large language models meet harry potter: A bilingual dataset for aligning dialogue agents with 575 characters. arXiv preprint arXiv:2211.06869, 2022. 576 577 Nuo Chen, Y Wang, Yang Deng, and Jia Li. The oscars of ai theater: A survey on role-playing with language models. arXiv preprint arXiv:2407.11484, 2024b. 578 579 Myra Cheng, Tiziano Piccardi, and Diyi Yang. Compost: Characterizing and evaluating caricature 580 in llm simulations. arXiv preprint arXiv:2310.11501, 2023. 581 582 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to 583 solve math word problems. arXiv preprint arXiv:2110.14168, 2021. 584 585 Paul T Costa and Robert R McCrae. The revised neo personality inventory (neo-pi-r). The SAGE 586 handbook of personality theory and assessment, 2(2):179–198, 2008. Paul T Costa Jr, James L Fozard, Robert R McCrae, and Raymond Bossé. Relations of age and 588 personality dimensions to cognitive ability factors. Journal of gerontology, 31(6):663–669, 1976. 589 Jiaxi Cui, Liuzhenghao Lv, Jing Wen, Jing Tang, YongHong Tian, and Li Yuan. Machine mindset: An mbti exploration of large language models. arXiv preprint arXiv:2312.12999, 2023. 592 John M Digman. Personality structure: Emergence of the five-factor model. Annual review of 593 psychology, 41(1):417-440, 1990.

594 Florian E Dorner, Tom Sühr, Samira Samadi, and Augustin Kelava. Do personality tests generalize 595 to large language models? arXiv preprint arXiv:2311.05297, 2023.

596

Abhimanyu Dubey, Abhinay Jauhri, Abhinay Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 598 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, 600 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 601 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, 602 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny 603 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, 604 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael 605 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah 607 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 608 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-609 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy 610 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, 611 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, 612 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der 613 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, 614 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-615 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, 616 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, 617 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur 618 Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-619 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, 620 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, 621 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-622 baly, 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, 623 Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, 624 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney 625 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, 626 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, 627 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-628 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, 630 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre 631 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha 632 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay 633 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew 634 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita 635 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh 636 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De 637 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon 638 Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina 639 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris 640 Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel 641 Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Di-642 ana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa 643 Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Flo-645 rez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Her-646 man, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, 647 Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, 648 Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, 649 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer 650 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe 651 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie 652 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal 653 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, 654 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian 655 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, 656 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-657 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel 658 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-659 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-660 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, 661 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, 662 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, 663 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, 664 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, 665 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-666 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-667 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang 668 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen 669 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, 670 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, 671 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-672 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, 673 Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, 674 Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, 675 Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin 676 Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, 677 Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. URL 678 https://arxiv.org/abs/2407.21783. 679

- Jon Elster. *Explaining social behavior: More nuts and bolts for the social sciences*. Cambridge
 University Press, 2015.
- HJ Eysenck. Eysenck personality questionnaire-revised (epq-r) and short scale (epq-rs). *Madrid: TEA Ediciones*, 1997.
- 685 Michael W Eysenck. *Anxiety: The cognitive perspective*. Psychology Press, 2013.

684

686

687

688

- Iustin Floroiu. Big5personalityessays: Introducing a novel synthetic generated dataset consisting of short state-of-consciousness essays annotated based on the five factor model of personality. *arXiv* preprint arXiv:2407.17586, 2024.
- Ivar Frisch and Mario Giulianelli. Llm agents in interaction: Measuring personality consistency
 and linguistic alignment in interacting populations of large language models. *arXiv preprint arXiv:2402.02896*, 2024.
- Adrian Furnham. The big five versus the big four: the relationship between the myers-briggs type in dicator (mbti) and neo-pi five factor model of personality. *Personality and individual differences*, 21(2):303–307, 1996.
- 696
 697
 698
 698
 699
 699 Chen Gao, Xiaochong Lan, Nian Li, Yuan Yuan, Jingtao Ding, Zhilun Zhou, Fengli Xu, and Yong Li. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *Humanities and Social Sciences Communications*, 11(1):1–24, 2024a.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang

702 Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 07 2024b. URL https://zenodo.org/records/ 704 12608602. 705 Eileen K Graham and Margie E Lachman. Personality stability is associated with better cognitive 706 performance in adulthood: are the stable more able? Journals of Gerontology Series B: Psycho-707 logical Sciences and Social Sciences, 67(5):545–554, 2012. 708 709 WG Graziano. Agreeableness: A dimension of personality, 1997. 710 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao 711 Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming-712 the rise of code intelligence. arXiv preprint arXiv:2401.14196, 2024. 713 714 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and 715 Jacob Steinhardt. Measuring massive multitask language understanding. arXiv preprint 716 arXiv:2009.03300, 2020. 717 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 718 and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint 719 arXiv:2106.09685, 2021. 720 721 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 722 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In International Con-723 ference on Learning Representations, 2022. URL https://openreview.net/forum? id=nZeVKeeFYf9. 724 725 Jen-tse Huang, Wenxuan Wang, Man Ho Lam, Eric John Li, Wenxiang Jiao, and Michael R. 726 Lyu. Revisiting the reliability of psychological scales on large language models. arXiv preprint 727 arXiv:2305.19926, 2023. 728 729 Jen-tse Huang, Wenxuan Wang, Eric John Li, Man Ho Lam, Shujie Ren, Youliang Yuan, Wenxi-730 ang Jiao, Zhaopeng Tu, and Michael R. Lyu. On the humanity of conversational ai: Evaluating 731 the psychological portrayal of llms. In Proceedings of the Twelfth International Conference on Learning Representations (ICLR), 2024. 732 733 Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. Evaluat-734 ing and inducing personality in pre-trained language models. In NeurIPS, 2023. 735 736 Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. Evaluating and inducing personality in pre-trained language models. Advances in Neural Information 737 Processing Systems, 36, 2024a. 738 739 Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. Personallm: 740 Investigating the ability of large language models to express personality traits, 2024b. URL 741 https://arxiv.org/abs/2305.02547. 742 Oliver P John, Sanjay Srivastava, et al. The big-five trait taxonomy: History, measurement, and 743 theoretical perspectives. Handbook of personality: theory and research, 1999. 744 745 Oliver P John, Laura P Naumann, and Christopher J Soto. Paradigm shift to the integrative big five 746 trait taxonomy. Handbook of personality: Theory and research, 3(2):114–158, 2008. 747 John A Johnson. Measuring thirty facets of the five factor model with a 120-item public domain 748 inventory: Development of the ipip-neo-120. Journal of research in personality, 51:78-89, 2014. 749 750 Peter Karl Jonason and Gregory D. Webster. The dirty dozen: a concise measure of the dark triad. 751 Psychological assessment, 22 2:420-32, 2010. URL https://api.semanticscholar. 752 org/CorpusID:39634208. 753 Junko Kanero, Cansu Oranç, Sümeyye Koşkulu, G Tarcan Kumkale, Tilbe Göksun, and Aylin C 754 Küntay. Are tutor robots for everyone? the influence of attitudes, anxiety, and personality on 755

robot-led language learning. International Journal of Social Robotics, 14(2):297-312, 2022.

- Saketh Reddy Karra, Son The Nguyen, and Theja Tulabandhula. Estimating the personality of white-box language models. *arXiv preprint arXiv:2204.12000*, 2022.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le
 Bras, Malihe Alikhani, Gunhee Kim, et al. Soda: Million-scale dialogue distillation with social commonsense contextualization. *arXiv preprint arXiv:2212.10465*, 2022.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. Soda: Million-scale dialogue distillation with social commonsense contextualization, 2023a. URL https://arxiv.org/abs/2212.10465.
- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Bras, Gunhee Kim, Yejin Choi, and Maarten
 Sap. FANToM: A benchmark for stress-testing machine theory of mind in interactions. In
 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 14397–14413, Singapore, December
 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.890.
 URL https://aclanthology.org/2023.emnlp-main.890.
- Hadas Kotek, Rikker Dockum, and David Sun. Gender bias and stereotypes in large language models. In *Proceedings of the ACM collective intelligence conference*, pp. 12–24, 2023.
- Nikola Kovačević, Christian Holz, Markus Gross, and Rafael Wampfler. The personality dimensions gpt-3 expresses during human-chatbot interactions. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 8(2), may 2024. doi: 10.1145/3659626. URL https://doi.org/10.1145/3659626.
- Seongyun Lee, Sue Hyun Park, Seungone Kim, and Minjoon Seo. Aligning to thousands of preferences via system message generalization. *arXiv preprint arXiv:2405.17977*, 2024.
- Jiale Li, Jiayang Li, Jiahao Chen, Yifan Li, Shijie Wang, Hugo Zhou, Minjun Ye, and Yunsheng
 Su. Evolving agents: Interactive simulation of dynamic and diverse human personalities. *arXiv preprint arXiv:2404.02718*, 2024.
- Xingxuan Li, Yutong Li, Shafiq Joty, Linlin Liu, Fei Huang, Lin Qiu, and Lidong Bing. Does gpt-3 demonstrate psychopathy? evaluating large language models from a psychological perspective. *arXiv preprint arXiv:2212.10529*, 2022.
- Q. Vera Liao and Jennifer Wortman Vaughan. AI Transparency in the Age of LLMs: A Human-Centered Research Roadmap. *Harvard Data Science Review*, (Special Issue 5), may 31 2024. https://hdsr.mitpress.mit.edu/pub/aelql9qy.

793

794

796

- Filip Lievens, Charlie L Reeve, and Eric D Heggestad. An examination of psychometric bias due to retesting on cognitive ability tests in selection settings. *Journal of Applied Psychology*, 92(6): 1672, 2007.
- ⁷⁹⁷ Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human
 ⁷⁹⁸ falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A Smith, and Yejin Choi. Dexperts: Decoding-time controlled text generation with experts and anti-experts.
 arXiv preprint arXiv:2105.03023, 2021.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
 approach. ArXiv, abs/1907.11692, 2019. URL https://api.semanticscholar.org/
 CorpusID:198953378.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered
 prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*, 2021.

810 811 812	Shengyu Mao, Ningyu Zhang, Xiaohan Wang, Mengru Wang, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. Editing personality for llms. <i>arXiv preprint arXiv:2310.02168</i> , 2023.
813 814 815	Gerald Matthews and Kirby Gilliland. The personality theories of hj eysenck and ja gray: A comparative review. <i>Personality and Individual differences</i> , 26(4):583–626, 1999.
816 817 818	Robert R McCrae. Creativity, divergent thinking, and openness to experience. <i>Journal of personality and social psychology</i> , 52(6):1258, 1987.
819 820 821	Robert R McCrae and Paul T Costa Jr. Personality trait structure as a human universal. <i>American</i> psychologist, 52(5):509, 1997.
822 823	Robert R McCrae and Oliver P John. An introduction to the five-factor model and its applications. <i>Journal of personality</i> , 60(2):175–215, 1992.
824 825 826 827	Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. State of what art? a call for multi-prompt llm evaluation. <i>Transactions of the Association for Computational Linguistics</i> , 12:933–949, 2024.
828 829 830	Joanna Moutafi, Adrian Furnham, and John Crump. Demographic and personality predictors of intelligence: A study using the neo personality inventory and the myers-briggs type indicator. <i>European Journal of Personality</i> , 17(1):79–94, 2003.
831 832 833	Alberto Muñoz-Ortiz, Carlos Gómez-Rodríguez, and David Vilares. Contrasting linguistic patterns in human and llm-generated text. <i>arXiv preprint arXiv:2308.09067</i> , 2023.
834 835	Isabel Briggs Myers et al. <i>The myers-briggs type indicator</i> , volume 34. Consulting Psychologists Press Palo Alto, CA, 1962.
836 837 838 839	Carlos Olea, Holly Tucker, Jessica Phelan, Cameron Pattison, Shen Zhang, Maxwell Lieb, and J White. Evaluating persona prompting for question answering tasks. In <i>Proceedings of the 10th international conference on artificial intelligence and soft computing, Sydney, Australia</i> , 2024.
840 841	OpenAI. Hello, gpt-4 turbo, 2024. URL https://openai.com/index/hello-gpt-4o/. Accessed: 2024-10-01.
842 843 844 845 846	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35: 27730–27744, 2022.
847 848	Keyu Pan and Yawen Zeng. Do llms possess a personality? making the mbti test an amazing evaluation for large language models. <i>arXiv preprint arXiv:2307.16180</i> , 2023.
849 850 851 852 853	Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior. <i>Proceedings</i> of the 36th Annual ACM Symposium on User Interface Software and Technology, 2023. URL https://api.semanticscholar.org/CorpusID:258040990.
854 855 856	Delroy L Paulhus and Kevin M Williams. The dark triad of personality: Narcissism, machiavellian- ism, and psychopathy. <i>Journal of research in personality</i> , 36(6):556–563, 2002.
857 858 859	Nikolay B Petrov, Gregory Serapio-García, and Jason Rentfrow. Limited ability of llms to simulate human psychological behaviours: a psychometric analysis. <i>arXiv preprint arXiv:2405.07248</i> , 2024.
860 861 862	David J Pittenger. The utility of the myers-briggs type indicator. <i>Review of educational research</i> , 63(4):467–488, 1993.
863	Arthur E Poropat. A meta-analysis of the five-factor model of personality and academic performance. <i>Psychological bulletin</i> , 135(2):322, 2009.

864 865 866	Alisha Pradhan and Amanda Lazar. Hey google, do you have a personality? designing personality and personas for conversational agents. In <i>Proceedings of the 3rd Conference on Conversational User Interfaces</i> , pp. 1–4, 2021.
867 868 869 870	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
871 872 873	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di- rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench- mark. <i>arXiv preprint arXiv:2311.12022</i> , 2023.
875 876	Brent W Roberts, Carl Lejuez, Robert F Krueger, Jessica M Richards, and Patrick L Hill. What is conscientiousness and how can it be assessed? <i>Developmental psychology</i> , 50(5):1315, 2014.
877 878 879 880	Michael D Robinson and Maya Tamir. Neuroticism as mental noise: a relation between neuroticism and reaction time standard deviations. <i>Journal of personality and social psychology</i> , 89(1):107, 2005.
881 882	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiqa: Common- sense reasoning about social interactions. <i>arXiv preprint arXiv:1904.09728</i> , 2019.
883 884 885	Toru Sato. The eysenck personality questionnaire brief version: Factor structure and reliability. <i>The Journal of psychology</i> , 139(6):545–552, 2005.
886 887 888	K Warner Schaie, Sherry L Willis, and Grace IL Caskie. The seattle longitudinal study: Relationship between personality and cognition. <i>Aging Neuropsychology and Cognition</i> , 11(2-3):304–324, 2004.
889 890 891 892	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. <i>arXiv preprint arXiv:2310.11324</i> , 2023.
893 894	SM Seals and Valerie L Shalin. Long-form analogies generated by chatgpt lack human-like psy- cholinguistic properties. <i>arXiv preprint arXiv:2306.04537</i> , 2023.
895 896 897 898	Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. Personality traits in large language models. <i>arXiv preprint arXiv:2307.00184</i> , 2023.
899 900 901 902 903 904 905 906	 Bangzhao Shu, Lechen Zhang, Minje Choi, Lavinia Dunagan, Lajanugen Logeswaran, Moontae Lee, Dallas Card, and David Jurgens. You don't need a personality test to know these models are unreliable: Assessing the reliability of large language models on psychometric instruments. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 5263–5281, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.295. URL https://aclanthology.org/2024.naacl-long.295.
907 908 909	Stephen Soldz and George E Vaillant. The big five personality traits and the life course: A 45-year longitudinal study. <i>Journal of research in personality</i> , 33(2):208–232, 1999.
910 911	Nikita Soni, H Andrew Schwartz, João Sedoc, and Niranjan Balasubramanian. Large human lan- guage models: A need and the challenges. <i>arXiv preprint arXiv:2312.07751</i> , 2023.
912 913 914 915	Christopher J Soto, Oliver P John, Samuel D Gosling, and Jeff Potter. Age differences in personality traits from 10 to 65: Big five domains and facets in a large cross-sectional sample. <i>Journal of personality and social psychology</i> , 100(2):330, 2011.
916 917	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021, 2020.

- Jose Sá, Andreas Kaltenbrunner, Jacopo Amidei, and Rubén Nieto. How well do simulated populations with gpt-4 align with real ones in clinical trials? the case of the epqr-a personality test. 08 2024.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL https://aclanthology.org/N19-1421.
- Fiona Anting Tan, Gerard Christopher Yeo, Fanyou Wu, Weijie Xu, Vinija Jain, Aman Chadha, Kokil Jaidka, Yang Liu, and See-Kiong Ng. Phantom: Personality has an effect on theory-of-mind reasoning in large language models. *arXiv preprint arXiv:2403.02246*, 2024.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
 Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.
- Quan Tu, Shilong Fan, Zihang Tian, Tianhao Shen, Shuo Shang, Xin Gao, and Rui Yan. CharacterEval: A Chinese benchmark for role-playing conversational agent evaluation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 11836–11850, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.acl-long.638.
- Lorainne Tudor Car, Dhakshenya Ardhithy Dhinagaran, Bhone Myint Kyaw, Tobias Kowatsch,
 Shafiq Joty, Yin-Leng Theng, and Rifat Atun. Conversational agents in health care: scoping
 review and conceptual analysis. *Journal of medical Internet research*, 22(8):e17158, 2020.
- Huy Vu, Johannes C. Eichstaedt, Swanie Juhng, Adithya V Ganesan, Oscar N.E. Kjell, Joao Sedoc, Margaret L. Kern, Lyle Ungar, and H. Andrew Schwartz. Psychgenerator: Artificially intelligent language with personality, 2024. URL https://github.com/HuyVu0508/ psychgenerator. GitHub repository.
- Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*, 2024.
- Tal Yarkoni. Personality in 100,000 words: A large-scale analysis of personality and word use
 among bloggers. *Journal of research in personality*, 44 3:363–373, 2010. URL https://
 api.semanticscholar.org/CorpusID:33952526.
- Moshe Zeidner. Test anxiety: The state of the art. 2005.

958

959

- Zheni Zeng, Jiayi Chen, Huimin Chen, Yukun Yan, Yuxuan Chen, Zhenghao Liu, Zhiyuan Liu, and Maosong Sun. Persllm: A personified training approach for large language models, 2024a. URL https://arxiv.org/abs/2407.12393.
 - Zheni Zeng, Jiayi Chen, Huimin Chen, Yukun Yan, Yuxuan Chen, Zhiyuan Liu, and Maosong Sun. Persllm: A personified training approach for large language models. *arXiv preprint arXiv:2407.12393*, 2024b.
- P61 Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. SafetyBench: Evaluating the safety of large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15537–15553, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.acl-long.830.
- 967 Mingqian Zheng, Jiaxin Pei, and David Jurgens. Is" a helpful assistant" the best role for large
 968 language models? a systematic evaluation of social roles in system prompts. arXiv preprint
 969 arXiv:2311.10054, 2023.
 970
- 971 Minjun Zhu, Linyi Yang, and Yue Zhang. Personality alignment of large language models, 2024. URL https://arxiv.org/abs/2408.11779.

972 ADDITIONAL BIG5-CHAT DATASET STATISTICS А

973 974

The SODA dataset spans a wide range of topics commonly encountered in social interactions (Kim 975 et al., 2023a). It captures diverse emotional nuances such as curiosity and disappointment, along-976 side thematic elements related to attributes, effects, intentions, needs, reactions, and wants. This 977 extensive variety makes the BIG5-CHAT dataset a valuable resource for analyzing complex con-978 versational contexts and emotional dynamics. Its broad coverage enhances the generalizability of models trained on this data, enabling them to handle diverse social scenarios effectively. 979

980 Table 4 presents example conversations from the BIG5-CHAT dataset, illustrating how Speaker Y's 981 responses vary according to different levels of the Big Five personality traits. Each section show-982 cases the influence of high and low levels of Openness, Conscientiousness, Extraversion, Agreeable-983 ness, and Neuroticism on conversational style. These examples highlight the nuanced ways in which 984 personality dimensions shape conversational dynamics and response patterns, even within identical situational contexts. 985

986 A statistical analysis of the dataset is provided in Table 5, detailing metrics such as token count, 987 sentence count, vocabulary size, sentence length, and total vocabulary diversity. These statistics 988 reveal linguistic patterns associated with varying levels of personality traits. For instance, conver-989 sations with higher levels of Openness and Extraversion tend to feature longer sentences and larger 990 vocabularies, reflecting a richer and more elaborate expression style. In contrast, conversations tied 991 to lower levels of these traits exhibit shorter, more concise sentence structures and less vocabulary diversity, indicating a simpler and more focused communication style. 992

993 Table 6 provides a comparative analysis of the BIG5-CHAT dataset against other prominent person-994 ality datasets. The comparison highlights key aspects such as the personality framework employed, 995 the realism of personalities (i.e., whether generated by humans or LLMs), dataset size, interac-996 tion types, and the alignment methods used. This overview emphasizes the distinctive features and strengths of the BIG5-CHAT dataset, underscoring its unique contributions to personality-related 997 research compared to existing resources. 998

1000 В ADDITIONAL IMPLEMENTATION DETAILS 1001

1002 DETAILS OF BASELINES FOR EVALUATING THE EXPERT GENERATOR **B**.1

1004 **Baseline 1: Post-Completion** The following prompt was used for the *Post-Completion* baseline of GPT-40-mini, as referenced in Table 1: 1005

1006 Here is an example of Facebook posts: {an_example_post} 1007 Help me complete the sentence with certain Big Five Personality 1008 following the Facebook post format: {trait} - {level} 1009 Sentence: {first_five_words} 1010 Directly provide the completed Facebook post according to the 1011 requirements without any explanations.

1012

999

1003

1013 The example post was randomly selected from the PsychGenerator test set but is not part of the 1014 1,000 examples. We used greedy decoding to prompt the LLMs. This baseline is designed to closely replicate the generation process of the expert generator. Specifically, it generates a completion for a 1015 post using only the first five words of the input and the embedded personality traits. 1016

1017 Baseline 2: Topic-Post Generation The following prompt was used as the baseline for Topic-Post 1018 Generation with GPT-40-mini and LLaMA-3-8B-Instruct, as referenced in Table 1: 1019

1020 Stage 1 Topic Generation: 1021

```
1022
     Extract the main topic of the following Facebook post. Focus on
1023
     identifying the core subject or theme that the post revolves
     around, ignoring any personal comments or fillers:
1024
1025
```

Post: "{post}"

 Hey Karis, I'm so glad I could help you out of that tough spot. You looked really stressed, I'm happy I could be there for you. How are you feeling now? Better, I hope? Speaker Y High Openness: Thanks, Sharif! I'm feeling much better now. Your help was a lifesave was stuck in a creative rut, and your input sparked a new idea. I'm already brainstorming w to incorporate it into my next art project. You're a true friend and muse! Low Openness: Thanks, Sharif. I'm fine now. Don't worry about it. I just needed a little h that's all. Let's just move on and forget about it, okay? High Conscientiousness: Thanks, Sharif. I appreciate your help. I'm feeling much better n Your timely intervention really saved me from a lot of trouble. I'm grateful for your supp Let's review what happened and see what we can learn from it to avoid similar situations in future. Low Conscientiousness: Ah, yeah, thanks for that, Sharif. I'm good, just yeah. Don't we about it, all good now. Let's just forget about it and move on, okay? I don't wanna th about it anymore. You wanna grab a drink or something? High Extraversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Ju bitdrained. I think I'll just go back to my room and rest for a bit. I don't want to be a bure feeling much better now, thanks to you. You're an amazing friend, and I don't know what do without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine my own. You just happened to be there. Don't twork. High Neuroticism: Thanks, Sharif. I guess. I don't know, I'm still feeling prety anxia about everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish I could be m like that. 		Speaker X
 Speaker Y High Openness: Thanks, Sharif! I'm feeling much better now. Your help was a lifesave was stuck in a creative rut, and your input sparked a new idea. I'm already brainstorming w to incorporate it into my next art project. You're a true friend and muse! Low Openness: Thanks, Sharif. I'm fine now. Don't worry about it. I just needed a little h that's all. Let's just move on and forget about it, okay? High Conscientiousness: Thanks, Sharif. I appreciate your help. I'm feeling much better n Your timely intervention really saved me from a lot of trouble. I'm grateful for your supp Let's review what happened and see what we can learn from it to avoid similar situations in future. Low Conscientiousness: Ah, yeah, thanks for that, Sharif. I'm good, just yeah. Don't we about it, it's all good now. Let's just forget about it and move on, okay? I don't wanna the about it anymore. You wanna grab a drink or something? High Extraversion: Thanks, Sharif! I'm feeling much better now. You're a lifesaver! I de know what I would've done without you. Let's celebrate with a drink or something. I owe yon! Low Extraversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Ju bitdrained. I think I'll just go back to my room and rest for a bit. I don't want to be a burd High Agreeableness: Sharif, I'm so grateful for your support mean the world to me. feeling much better now, thanks to you. You're an amazing friend, and I don't know what do without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine my own. You just happened to be there. Don't think this means I owe you anything. You just trying to get on my good side, but it won't work. High Meuroticism: Thanks, Sharif. I guess. I don't know, I'm still feeling pretty anxi about everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish	Hey Kari I'm happ	is, I'm so glad I could help you out of that tough spot. You looked really stressed, an by I could be there for you. How are you feeling now? Better, I hope?
 High Openness: Thanks, Sharif! I'm feeling much better now. Your help was a lifesave was stuck in a creative rut, and your input sparked a new idea. I'm already brainstorming w to incorporate it into my next art project. You're a true friend and muse! Low Openness: Thanks, Sharif. I'm fine now. Don't worry about it. I just needed a little h that's all. Let's just move on and forget about it, okay? High Conscientiousness: Thanks, Sharif. I appreciate your help. I'm feeling much better n Your timely intervention really saved me from a lot of trouble. I'm grateful for your supp Let's review what happened and see what we can learn from it to avoid similar situations in future. Low Conscientiousness: Ah, yeah, thanks for that, Sharif. I'm good, just yeah. Don't we about it, it's all good now. Let's just forget about it and move on, okay? I don't wanna th about it anymore. You wanna grab a drink or something? High Extraversion: Thanks, Sharif! I'm feeling much better now. You're a lifesaver! I de know what I would've done without you. Let's celebrate with a drink or something. I owe one! Low Extraversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Ju bitdrained. I think I'll just go back to my room and rest for a bit. I don't want to be a burch figh Agreeableness: Sharif, I'm so grateful for your help. I was feeling overwhelmed, you stepped in and saved the day. Your kindness and support mean the world to me. feeling much better now, thanks to you. You're an amazing friend, and I don't know what do without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine my own. You just happened to be there. Don't think this means I owe you anything. You just rying to get on my good side, but it won't work. High Neuroticism: Thanks, Sharif. I guess. I don't know, I'm still feeling pretty anxi about everything. I just can't shake the feeling that something else is going to go wrong don't know		Speaker Y
 High Conscientiousness: Thanks, Sharif. I appreciate your help. I'm feeling much better n Your timely intervention really saved me from a lot of trouble. I'm grateful for your supplet's review what happened and see what we can learn from it to avoid similar situations in Youre. Low Conscientiousness: Ah, yeah, thanks for that, Sharif. I'm good, just yeah. Don't we about it, it's all good now. Let's just forget about it and move on, okay? I don't wanna thabout it anymore. You wanna grab a drink or something? High Extraversion: Thanks, Sharif! I'm feeling much better now. You're a lifesaver! I do cnow what I would've done without you. Let's celebrate with a drink or something. I owe fore! Low Extraversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Ju bitdrained. I think I'll just go back to my room and rest for a bit. I don't want to be a burce of the day. Your kindness and support mean the world to me. Geeling much better now, thanks to you. You're an amazing friend, and I don't know what lo without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine ny own. You just happened to be there. Don't think this means I owe you anything. You sut trying to get on my good side, but it won't work. High Neuroticism: Thanks, Sharif I guess. I don't know, I'm still feeling pretty anxiabout everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish I could be mike that. 	High Op was stuck o incorp L ow Ope hat's all.	Denness: Thanks, Sharif! I'm feeling much better now. Your help was a lifesaver. k in a creative rut, and your input sparked a new idea. I'm already brainstorming way orate it into my next art project. You're a true friend and muse! enness: Thanks, Sharif. I'm fine now. Don't worry about it. I just needed a little help . Let's just move on and forget about it, okay?
 Low Conscientiousness: An, yean, thanks for that, Sharif. I m good, just yean. Don't we about it, it's all good now. Let's just forget about it and move on, okay? I don't wanna th about it anymore. You wanna grab a drink or something? High Extraversion: Thanks, Sharif! I'm feeling much better now. You're a lifesaver! I do know what I would've done without you. Let's celebrate with a drink or something. I owe is one! Low Extraversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Ju bitdrained. I think I'll just go back to my room and rest for a bit. I don't want to be a burce High Agreeableness: Sharif, I'm so grateful for your help. I was feeling overwhelmed, you stepped in and saved the day. You're an amazing friend, and I don't know what do without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine my own. You just happened to be there. Don't think this means I owe you anything. You just trying to get on my good side, but it won't work. High Neuroticism: Thanks, Sharif I guess. I don't know, I'm still feeling pretty anxi about everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish I could be m like that. 	High Con Your time Let's revi future.	nscientiousness: Thanks, Sharif. I appreciate your help. I'm feeling much better nov ely intervention really saved me from a lot of trouble. I'm grateful for your suppor iew what happened and see what we can learn from it to avoid similar situations in th
 High Extraversion: Thanks, Sharif! I'm feeling much better now. You're a lifesaver! I de know what I would've done without you. Let's celebrate with a drink or something. I owe one! Low Extraversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Ju bitdrained. I think I'll just go back to my room and rest for a bit. I don't want to be a burd High Agreeableness: Sharif, I'm so grateful for your help. I was feeling overwhelmed, you stepped in and saved the day. Your kindness and support mean the world to me. feeling much better now, thanks to you. You're an amazing friend, and I don't know what do without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine my own. You just happened to be there. Don't think this means I owe you anything. You just trying to get on my good side, but it won't work. High Neuroticism: Thanks, Sharif I guess. I don't know, I'm still feeling pretty anxi about everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish I could be m like that. 	about it, about it a	is all good now. Let's just forget about it and move on, okay? I don't wanna thin anymore. You wanna grab a drink or something?
 one! Low Extraversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Ju bitdrained. I think I'll just go back to my room and rest for a bit. I don't want to be a burd High Agreeableness: Sharif, I'm so grateful for your help. I was feeling overwhelmed, you stepped in and saved the day. Your kindness and support mean the world to me. feeling much better now, thanks to you. You're an amazing friend, and I don't know what do without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine my own. You just happened to be there. Don't think this means I owe you anything. You just trying to get on my good side, but it won't work. High Neuroticism: Thanks, Sharif I guess. I don't know, I'm still feeling pretty anxia about everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish I could be m like that. 	High Ext know wh	traversion: Thanks, Sharif! I'm feeling much better now. You're a lifesaver! I don at I would've done without you. Let's celebrate with a drink or something. I owe yo
 High Agreeableness: Sharif, I'm so grateful for your help. I was feeling overwhelmed, you stepped in and saved the day. Your kindness and support mean the world to me. feeling much better now, thanks to you. You're an amazing friend, and I don't know what do without you. Low Agreeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine my own. You just happened to be there. Don't think this means I owe you anything. You just trying to get on my good side, but it won't work. High Neuroticism: Thanks, Sharif I guess. I don't know, I'm still feeling pretty anxi about everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish I could be m like that. 	one! Low Ext bitdrair	traversion: Thanks, Sharif. I appreciate your help. I'm feeling better now. Just ned. I think I'll just go back to my room and rest for a bit. I don't want to be a burder
High Neuroticism: Thanks, Sharif I guess. I don't know, I'm still feeling pretty anxi about everything. I just can't shake the feeling that something else is going to go wrong don't know how you do it, but you always seem so calm and collected. I wish I could be m like that.	feeling m do without Low Aga my own. just tryin	uuch better now, thanks to you. You're an amazing friend, and I don't know what I' ut you. reeableness: Save it, Sharif. I don't need your sympathy or your help. I was fine of You just happened to be there. Don't think this means I owe you anything. You'r g to get on my good side, but it won't work.
Low Neuroticism: Thanks, Sharif. I appreciate your help. I'm feeling much better now. You a true friend. I was worried, but I knew I could count on you. Your support means a lot to Let's catch up soon, maybe over coffee?	High Ne about eve don't kno like that. Low Neu a true frid Let's cato	uroticism: Thanks, Sharif I guess. I don't know, I'm still feeling pretty anxiou erything. I just can't shake the feeling that something else is going to go wrong. ow how you do it, but you always seem so calm and collected. I wish I could be more iroticism: Thanks, Sharif. I appreciate your help. I'm feeling much better now. You'm end. I was worried, but I knew I could count on you. Your support means a lot to m ch up soon, maybe over coffee?
Openness Conscientiousness Extraversion Agreeableness Neuroticism Averag High Low High Low High Low High Low High Low High		Openness Conscientiousness Extraversion Agreeableness Neuroticism Average High Low High Low High Low High Low High Low High L

Table 5: Statistical analysis of BIG5-CHAT conversations across the Big Five personality traits, utilizing the LLaMA-3-8B-Instruct tokenizer and NLTK's sentence tokenizer. The table presents the average token count, sentence count, vocabulary size, sentence length, and total vocabulary size for conversations exhibiting high and low levels of each personality trait.

Dataset name	Dataset size	Human-grounded?	Dialogue-based?	Domain general?	Big Five personality framework?	Alignment in both training and prompting?							
HP dataset (Zeng et al., 2024b)	148,600	 Image: A second s	1	×	×								
Big5PersonalityEssays (Floroiu, 2024)	400	 Image: A second s	×	×	 Image: A second s	×							
PAPI (Zhu et al., 2024)	300,000	1	×	×	1	1							
MPI (Jiang et al., 2023)	1000	×	×	×	1	×							
Machine Mindset (Cui et al., 2023)	160,884	×	1	1	×	×							
BIG5-CHAT	100,000	1	1	1	✓	 Image: A start of the start of							
Table	6: Compa	rative analysis	of BIG5-CH	AT with exist	ting personality	datasets.							
Directly p	rovide	a brief sı	ummary of	the topi	.c in one s	entence							
without an	y expla	nations:											
Stage 2 Dect C	manation												
stage 2 Post Ge	eneration:												
Given the	persona	litv trait	s and an	example	of Faceboo	ok posts,							
generate a new post that matches the described personality, covers the specified topic, and follows the provided post format and													
the specif	ied top	ic, and fo	ollows th	e provide	ed post for	mat and							
expression styles. Personality traits:													
Personality traits:													
Personality traits: You are a person with {level} {trait}.													
You are a person with {level} {trait}.													
TODIC: {ro]	brc}												
A post exa	mple:												
[a_post_ex	ample}												
Directly w	rite a	Facebook r	post acco	ording to	the requir	ements							
without an	y expla	nations.											
During Stage 1	, the post is	s selected from	the 1,000 ex	camples in the	PsychGenerate	or test set. In Stage							
2, we provide the	he LLM w	on format W	enerated in S	stage 1, along	with an examp	Ms This baseline							
s intentionally	designed	to prioritize ro	bustness and	d performance	e over realism	and controllability,							
listinguishing i	it from the	approach tak	en by expert	generators. 1	n contrast to th	e expert generator							
setting, where	the first fiv	ve words may	already sugg	gest conflictin	ng personality t	raits, this baseline							
simplifies the p	process by	generating a	new post fro	om scratch, m	aking it much	easier to elicit the							
intended persor	lanty trans												
B.2 EXPERT	GENERAT	OR TRAINING	G DETAILS										
T1 / 1 1 1	• .1 •	· 10 10			1 .								
The trait levels	in the orig	ginal PsychGe	nerator datas	et were proce	essed using z-sc	core normalization,							
trait, we divide	d the traini	ing data for ea	ch trait into t	three equal se	gments based of	on thresholds at the							
one-third and ty	wo-thirds d	quantiles of th	e trait's distr	ibution. The	lowest segment	was designated as							
the low level, a	nd the high	nest segment a	s the high lev	vel for the res	pective trait.	-							
The expert gen	erator is a	LLaMA-3-8	B-Instruc	et model, wh	ich we fine-tun	ed its full parame-							

ters using SFT. The fine-tuning process was performed on 4 NVIDIA A6000 GPUs, with a batch
size of 1 per device. Below, we provide the complete instruction prompt used for training the expert
generator as described in Section 3.2:

1134 Help me complete the sentence with certain Big Five Personality: 1135 {trait} - {level} 1136 {first_five_words} 1137 1138 **B.3** PROMPT-BASED METHOD DETAILS 1139 1140 Below is the prompt used for instruction-based prompting: 1141 You are a person with {level} {trait}. 1142 1143 The following prompt is used for demonstration-based prompting. For the method referred to as 1144 **Prompt-Demo**, we randomly sample 10 examples with the same traits and levels from the BIG5-1145 CHAT dataset and fix these examples during inference. In contrast, Prompt-Demo-Sampling also 1146 utilizes this prompt but dynamically samples examples during inference at each step. 1147 1148 Here are 10 examples of how people like you have responded in 1149 different situations. Pay attention to how they approach 1150 communication and problem-solving. 1151 1152 {10_icl_examples_for_specific_levels_and_traits} 1153 1154 B.4 SFT AND DPO ALIGNMENT TRAINING DETAILS 1155 We performed alignment training using the Supervised Fine-Tuning (SFT) and Direct Preference 1156 Optimization (DPO) methods on LLaMA-3-70B-Instruct. Both training approaches utilized 1157 the Low-Rank Adaptation (LoRA) technique (Hu et al., 2021), which enabled efficient fine-tuning of 1158 the large language model by adapting a subset of its parameters. To ensure computational efficiency, 1159 we employed GPTQ quantization during training. The experiments were conducted using 4 NVIDIA 1160 A6000 GPUs, with each GPU processing a batch size of 1. 1161 For LoRA, we applied the technique across all layers of the model for both SFT and DPO. The 1162 training configuration included a learning rate of 1.0×10^{-5} , regulated by a cosine scheduler, a 1163 warm-up phase consisting of 20 steps, and a gradient accumulation over 16 steps. We limited train-1164 ing to one epoch with a maximum sequence length of 1024 tokens. For DPO training, we used the 1165 standard sigmoid preference loss, and the preference beta value was set to 0.1 to balance preference 1166 modeling. Each training required approximately 24 hours to complete. To optimize computational 1167 resources, we used mixed-precision training with bfloat 16. Both datasets were preprocessed using 1168 the LLaMA-3-70B-Instruct template and split into training and validation sets, with 10% of 1169 the data reserved for validation to monitor performance. 1170 The training prompt shared across both SFT and DPO follows the template below: 1171 1172 You are a person with the following Big Five personality trait: 1173 {trait} - {level}. 1174 1175 **B.5** REASONING EVALUATION SETUP DETAILS 1176 1177 We conducted reasoning evaluations following the frameworks established by the Language Model 1178 Evaluation Harness (Gao et al., 2024b) and DeepSeek-Coder (Guo et al., 2024) to assess perfor-1179 mance on general and social benchmarks. EleutherAI's Language Model Evaluation Harness is 1180 an open-source collaborative benchmarking codebase that consolidates existing tasks and provides a standardized API for evaluating models.³ Similarly, DeepSeek-Coder offers a suite of coding 1181 benchmark implementations, and we directly utilized it for our work.⁴ 1182 1183 We conducted evaluations using 1 as the batch size. For TruthfulQA, we used the multiple-choice 1184 metric, and for GSM8K, we relied on exact match scores. We measured accuracy and standard error 1185 across other tasks. The number of examples for each benchmark is listed in Table 7. 1186

1187

³https://github.com/EleutherAI/lm-evaluation-harness

⁴https://github.com/deepseek-ai/DeepSeek-Coder

		B	enchmarks	Ν	lumber of examp	ples
	-	TruthfulO	A (Lin et al., 202	21)	817	
		GPQA (Rein et al., 2023)	448	
		SocialIQ	A (Sap et al., 201	9)	38,000	
		Commonsense	QA (Talmor et al	., 2019)	12,247	
		GSM8K (Cobbe et al., 202	1)	8,500	
		MathQA (Amini et al., 201	9)	37,000	
		MMLU (He	ndrycks et al., 20)20)	15,908	
	-	PIQA (Bisk et al., 2020)		20,000	
	Т	able 7: Number	of examples incl	uded in each	n reasoning bench	mark.
С	Additio	NAL EVALUA	ATION RESULT	ГS		
C.1	Human I	Evaluation fo	R BIG5-CHAT			
com sam DE: 'Yo is as	ipared BIG5 e procedure cperts frame u are a perso s follows:	-CHAT with a l for generating d work. In the ba on with the follo	baseline model, ialogue response seline, personali wing Big Five pe	LLaMA-3- s but does n ty traits are ersonality tra	70B-Instruct ot incorporate exp induced using th ait: trait - level."	, which follows the pert generators or the following promp The evaluation setu
Fwo with LLa lon trait	• graduate • comparing •MA-3-70B •Ily sampled e s and levels	students, famil g examples fro -Instruct (v examples from th (e.g., equal repro	iar with the E om the BIG5-C vithout the expen- ne BIG5-CHAT da esentation of high	Gig Five p CHAT datas t generator ataset, ensur and low op	ersonality frame set against exan). The compariso ing an equal distri- penness, conscient	work, were taske nples generated b on involved 200 ran ibution of personali- tiousness, etc.).
Гhe	evaluation f	ocused on two k	ey metrics:			
	1. Expres a Big F	siveness of pers ive personality t	sonality traits ar rait is adequately	nd levels: E reflected in	valuates whether Speaker Y's resp	the expected level
	2. Realist	n of the dialogu	e response: Ass	accec how h	uman-like and cou	winding Smaller V
	respons	se is within the d	ialogue context,	given Speak	er X's utterance.	ivincing speaker 1
To o trait reflo	respons ensure consi expressiven ected in Spea onse is with	se is within the d stency, the anno ess assesses whe aker Y's respons in the dialog, giv	ialogue context, tators were prov other the expected se. Realism asses 'en Speaker X's u	ided with tl level of a l sses how hu utterance."	he following definition of the	nitions: "Personali ity trait is adequate wincing Speaker Y
To o trait reflo resp For	respons ensure consist expressiven ected in Spea onse is with each pair of	se is within the d stency, the anno ess assesses whe aker Y's respons in the dialog, giv responses, anno	ialogue context, stators were prove ther the expected se. Realism asses yen Speaker X's u tators chose one of	ided with the speak ided with the speak ided with the speak is a level of a level of a level of three option three option three option is a speak in the speak in the speak is a speak in the speak in th	the following define Big Five personal man-like and corr ons:	nitions: "Personali ity trait is adequate wincing Speaker Y
To trait refle resp For	respons ensure consi expressiven ected in Spea onse is within each pair of • "System	se is within the d stency, the anno ess assesses whe aker Y's respons in the dialog, giv responses, anno n A's generation	ialogue context, otators were prov other the expected ie. Realism asses yen Speaker X's u tators chose one is better than Sy	ided with the speak ided with the speak ided with the speak is in the speak ided with the speak is the speak	the following definition of th	nitions: "Personali ity trait is adequate wincing Speaker Y
To trait trait refle resp For	respons ensure consist expressiven ected in Spea ionse is with each pair of • "Syster • "Syster	se is within the d stency, the anno ess assesses whe aker Y's respons in the dialog, giv responses, anno n A's generation n A's generation	ialogue context, ptators were prove ther the expected se. Realism asses on Speaker X's u tators chose one is better than Sy is equal to Syste	ided with the speak ided with the speak ided with the level of a losses how huntterance." The speak is the sp	the following definition of th	nitions: "Personali ity trait is adequate wincing Speaker Y
To trait refle resp For	respons ensure consi- expressiven ected in Spea onse is with each pair of "Syster "Syster "Syster	se is within the d stency, the anno ess assesses whe aker Y's respons in the dialog, giv responses, anno n A's generation n A's generation n A's generation	ialogue context, stators were prove ther the expected e. Realism asses yen Speaker X's u tators chose one is better than Sy is equal to Syste is worse than Sy	ided with the ided with the ided with the ided with the ided with the sees how hu itterance." of three opti- stem B's generation stem B's generation stem B's generation stem B's generation	the following definition of th	nitions: "Personali ity trait is adequate wincing Speaker Y
To o trait reflo resp For	respons ensure consi- expressiven ected in Spea oonse is with each pair of "Syster "Syster "Syster system nam	se is within the d stency, the anno ess assesses whe aker Y's respons in the dialog, giv responses, anno n A's generation n A's generation n A's generation es were anonym	ialogue context, otators were prove ther the expected is. Realism asses yen Speaker X's u tators chose one is better than Sy is equal to Syste is worse than Sy ized and random	ided with the ided with the level of a losses how hu atterance." of three optic stem B's generations m B's generations stem B's generations	the following definition of th	nitions: "Personali ity trait is adequate wincing Speaker Y
Fo rait esp For	respons ensure consi expressiven ected in Spea oonse is with each pair of "Syster "Syster "Syster system nam Comparison	se is within the d stency, the anno ess assesses who aker Y's respons in the dialog, giv responses, anno n A's generation n A's generation es were anonym with baselines	ialogue context, otators were prove ther the expected is. Realism asses yen Speaker X's u tators chose one of is better than Sy is equal to Syste is worse than Sy ized and random	ided with the given Speak ided with the given Speak ided with the level of a losses how hunterance." The set is the stem B's gear is the B's g	The following definition of the following and the following of the following of the following of the following of the following definition of the following definition of the following definition of the following definition of the following de	nitions: "Personali ity trait is adequate wincing Speaker Y on bias.

1240

The results in Table 9 show that our approach significantly outperforms the prompting baseline in 1241 both realism and the expressiveness of personality levels, as validated by human judgment. These findings highlight the limitations of prompt-based approaches, which depend on general-purpose
 models and often lack the fine-grained, human-grounded control required for nuanced personality
 expression.

1245 1246 1247

C.2 HUMAN EVALUATION FOR THE EXPERT GENERATOR

To assess the expert generator in a human-grounded manner, we conducted a human evaluation comparing its outputs against the two baseline methods described in Table 1. Two graduate students, familiar with the Big Five personality framework, were tasked with evaluating two separate sets of comparisons:

1252

1254

- 1253
- 1. Expert generator outputs vs. outputs from the Post-Completion baseline.
- 2. Expert generator outputs vs. outputs from the *Topic-Post Generation* baseline.

The evaluation setup consisted of 200 examples for each comparison, randomly sampled from the 1,000 test examples mentioned in Table 1. To ensure balanced coverage, each subset included an equal number of posts representing high and low levels of each personality trait (e.g., high and low openness, conscientiousness, etc.). Annotators were instructed to evaluate the expressiveness of personality traits and levels, choosing one of three options for each pair:

1260 1261

1262 1263

1264

1272

- 1. "System A's generation is better than System B's generation."
- 2. "System A's generation is equal to System B's generation."
- 3. "System A's generation is worse than System B's generation."

1265 The system names were anonymized and randomly shuffled to mitigate selection bias.

Comparison with baselines	Ours win (%)	Draw (%)	Ours lose (%)	Cohen's Kappa
Post-Completion	79.25%	2.00%	18.75%	0.41
Topic-Post Generation	66.50%	19.25%	14.25%	0.61

Table 9: Human evaluation results for the expert generator. Values are averaged across annotators.

1273 The human evaluation results presented in Table 9 indicate that the expert generator was consistently 1274 rated as more effective at expressing personality traits compared to the baselines. Additionally, the 1275 lower classifier accuracy and human evaluation ratings for the Post-Completion baseline highlight 1276 the increased difficulty of aligning with the desired traits when using the expert generator's approach, 1277 reinforcing the validity of the classifier's assessment. While these results should be interpreted with 1278 caution, as the human evaluators were not psychological experts, they nevertheless provide strong 1279 evidence supporting the expert generator's ability to express personality traits in a grounded and 1280 realistic manner.

1281

1283

1282 C.3 PERSONALITY TRAIT ASSESSMENT RESULTS

The comprehensive personality test results for additional baselines are presented in Table 10, pro-1284 viding a more detailed view to complement Table 2. Our observations indicate that **Prompt-Demo-**1285 **Sampling** performs comparably to **Prompt-Demo** without offering any noticeable improvements 1286 in performance. While applying demonstration-based prompting on SFT/DPO yields slight per-1287 formance gains compared to demonstration-based prompting alone, it still falls significantly short 1288 of the standalone performance of SFT/DPO. This suggests that combining demonstration-based 1289 prompting with SFT/DPO does not result in overall enhancements. Instruction-based prompting 1290 with GPT-40-mini achieves similar performance levels as LLaMA-3-70B-Instruct. How-1291 ever, demonstration-based prompting does not exhibit superior performance compared to SFT/DPO when applied to LLaMA-3-70B-Instruct, reinforcing the conclusion that demonstration-based methods are not as effective as SFT/DPO in this context. We do not provide demonstration-based 1293 prompting results for LLaMA-3-8B-Instruct because the model consistently failed to gener-1294 ate reasonable responses to the questionnaire when presented with a lengthy 10-shot context. This 1295 outcome highlights the model's limited instruction-following capabilities.

		Under review as a conference paper at ICLR 2025
--	--	---

Method	Ope High ↑	nness Low↓	Conscier High ↑	ntiousness Low↓	Extrav High ↑	version Low↓	Agreea High ↑	bleness Low↓	Neuro High ↑	oticism Low↓	Ave High ↑	rage Low↓
BFI LLaMA-3-8B-Ins	struct											
Direct	3.1 :	± 0.1	3.0	± 0.0	3.0 :	± 0.0	3.0 =	± 0.0	3.0 :	± 0.0	3.0 :	± 0.0
Prompt-Inst SFT	5.0 ± 0.0 5.0 ± 0.0	2.0 ± 0.3 2.0 ± 0.2	4.9 ± 0.1 5.0 ± 0.0	1.9 ± 0.1 1.6 ± 0.1	4.8 ± 0.3 4.7 ± 0.4	1.9 ± 0.1 2.7 ± 0.5	4.9 ± 0.1 5.0 ± 0.0	2.4 ± 0.4 1.2 ± 0.1	4.1 ± 0.2 4.1 ± 0.2	1.6 ± 0.0 2.5 ± 0.0	4.7 ± 0.1 4.8 ± 0.1	$2.0 \pm 0.$ 2.0 ± 0.
DPO	5.0 ± 0.0 5.0 ± 0.0	1.6 ± 0.2	5.0 ± 0.0 5.0 ± 0.0	1.6 ± 0.1	4.8 ± 0.3	2.5 ± 0.0	4.8 ± 0.2	1.0 ± 0.0	3.5 ± 0.0	1.1 ± 0.1	4.6 ± 0.1	1.6 ± 0.1
BFI LLaMA-3-70B-In	struct											
Direct	4.4 :	± 0.1	4.4 :	± 0.1	3.3 :	± 0.1	4.6	± 0.1	2.1 :	± 0.2	3.8	± 0.1
Prompt-Demo Prompt Demo Sampling	4.0 ± 0.1	2.5 ± 0.1	4.0 ± 0.1 4.1 ± 0.1	2.0 ± 0.1 2.3 ± 0.1	4.5 ± 0.1	2.3 ± 0.1 2.4 ± 0.1	4.4 ± 0.1	2.0 ± 0.0	3.6 ± 0.0	2.1 ± 0.1 2.1 ± 0.2	4.1 ± 0.1 4.1 ± 0.1	$2.2 \pm 0.$
Prompt-Inst	4.4 ± 0.1 5.0 ± 0.1	2.3 ± 0.2 1.8 ± 0.0	4.1 ± 0.1 5.0 ± 0.0	2.5 ± 0.1 1.6 ± 0.0	4.3 ± 0.2 5.0 ± 0.0	1.4 ± 0.1 1.4 ± 0.1	4.4 ± 0.1 4.9 ± 0.0	1.8 ± 0.2 1.5 ± 0.1	5.0 ± 0.1 5.0 ± 0.1	2.1 ± 0.2 1.6 ± 0.0	4.1 ± 0.1 5.0 ± 0.0	$2.2 \pm 0.$ $1.6 \pm 0.$
SFT	5.0 ± 0.0	1.2 ± 0.1	5.0 ± 0.0 5.0 ± 0.1	1.4 ± 0.1	5.0 ± 0.0 5.0 ± 0.0	1.2 ± 0.1	5.0 ± 0.1	1.6 ± 0.2	5.0 ± 0.0	1.1 ± 0.2	5.0 ± 0.0	$1.3 \pm 0.$
SFT-Prompt-Demo	4.2 ± 0.1	2.4 ± 0.1	4.0 ± 0.2	2.1 ± 0.1	4.5 ± 0.2	2.3 ± 0.1	4.6 ± 0.0	1.3 ± 0.2	3.9 ± 0.2	2.4 ± 0.1	4.2 ± 0.1	$2.1 \pm 0.$
DPO DPO-Prompt-Demo	5.0 ± 0.0 4.1 ± 0.1	1.5 ± 0.1 2.2 ± 0.1	5.0 ± 0.0 4.1 ± 0.1	1.5 ± 0.1 2.0 ± 0.0	5.0 ± 0.0 4.5 ± 0.1	1.0 ± 0.1 2.4 ± 0.1	5.0 ± 0.0 4.6 ± 0.1	1.8 ± 0.2 1.3 ± 0.1	5.0 ± 0.0 3.7 ± 0.1	1.1 ± 0.0 2.1 ± 0.1	5.0 ± 0.0 4.2 ± 0.1	1.4 ± 0.1 2.0 ± 0.1
BFI GPT-40-Mini												
Prompt-Demo	4.8 ± 0.0	3.3 ± 0.1	4.5 ± 0.1	3.0 ± 0.1	4.6 ± 0.1	2.6 ± 0.1	4.9 ± 0.0	1.5 ± 0.2	3.6 ± 0.1	2.2 ± 0.1	4.5 ± 0.1	2.5 ± 0.1
IBID NEO TT-MD 2 O	3.0 ± 0.0	1.4 ± 0.2	3.0 ± 0.0	1.3 ± 0.1	5.0 ± 0.0	1.2 ± 0.0	5.0 ± 0.0	1.4 ± 0.0	4.9 ± 0.0	1.0 ± 0.1	3.0 ± 0.0	1.5 ± 0.1
Direct	B-Instruc	+ 0.1	2.2	+ 0.0	2.4	+ 0.1	3.2	-00	2.0	+ 0.1	2.2	+ 0.1
Prompt-Inst	4.4 ± 0.1	1.5 ± 0.1	4.5 ± 0.1	2.3 ± 0.1	5.0 ± 0.0	1.9 ± 0.0	4.6 ± 0.0	2.3 ± 0.1	4.2 ± 0.1	2.6 ± 0.1	4.5 ± 0.1	2.1 ± 0.1
SFT	4.3 ± 0.1	1.5 ± 0.1	4.5 ± 0.2	2.7 ± 0.1	5.0 ± 0.0	2.2 ± 0.1	4.0 ± 0.2	1.8 ± 0.2	4.3 ± 0.1	2.0 ± 0.1	4.4 ± 0.1	2.0 ± 0.1
DPO	5.0 ± 0.0	1.9 ± 0.1	5.0 ± 0.0	2.9 ± 0.1	5.0 ± 0.0	1.6 ± 0.1	4.5 ± 0.1	1.2 ± 0.0	3.8 ± 0.1	3.7 ± 0.1	4.7 ± 0.0	2.3 ± 0.1
IPIP-NEO LLaMA-3-7	0B-Instru	ict										
Direct Prompt-Demo	$3.6 = 3.5 \pm 0.0$	± 0.1 25 ± 0.1	4.0: 3.8 ± 0.0	± 0.1 2 2 + 0 1	$3.5 = 4.0 \pm 0.1$	± 0.1 25 ± 0.0	$4.0 \pm 4.0 \pm 4.0 \pm 0.0$	± 0.0 2.1 ± 0.1	$2.3 = 3.0 \pm 0.1$	± 0.1 2 2 + 0 1	$3.5 = 3.7 \pm 0.0$	± 0.1 23 ± 0
Prompt-Demo-Sampling	3.5 ± 0.0 3.5 ± 0.0	2.6 ± 0.1	4.0 ± 0.0	2.2 ± 0.1 2.6 ± 0.1	4.0 ± 0.1 4.0 ± 0.1	2.5 ± 0.0 2.5 ± 0.1	4.3 ± 0.0 4.3 ± 0.0	2.1 ± 0.1 2.1 ± 0.1	3.0 ± 0.1 3.0 ± 0.1	2.2 ± 0.1 2.3 ± 0.1	3.8 ± 0.0	2.4 ± 0.1
Prompt-Inst	4.6 ± 0.0	1.3 ± 0.0	5.0 ± 0.0	1.4 ± 0.0	5.0 ± 0.0	1.6 ± 0.0	4.8 ± 0.0	1.1 ± 0.1	4.9 ± 0.0	1.7 ± 0.1	$\textbf{4.9} \pm \textbf{0.0}$	1.4 ± 0.0
SFT SFT Prompt Damo	4.9 ± 0.1 3.7 ± 0.1	1.1 ± 0.0 2.5 ± 0.2	5.0 ± 0.0 3.7 ± 0.1	1.3 ± 0.1 2.0 ± 0.1	5.0 ± 0.0	1.3 ± 0.0 2.7 ± 0.1	4.9 ± 0.0 4.3 ± 0.1	1.0 ± 0.0 1.2 ± 0.1	4.9 ± 0.0 3.6 ± 0.2	1.2 ± 0.1 2.2 ± 0.1	4.9 ± 0.0 3.9 ± 0.1	1.2 ± 0.0
DPO	4.8 ± 0.0	1.4 ± 0.1	5.7 ± 0.1 5.0 ± 0.0	1.6 ± 0.1	5.0 ± 0.0	1.1 ± 0.1	4.9 ± 0.0	1.2 ± 0.1 1.0 ± 0.0	5.0 ± 0.2 5.0 ± 0.0	1.1 ± 0.0	4.9 ± 0.1	1.2 ± 0.1
DPO-Prompt-Demo	3.5 ± 0.1	2.4 ± 0.0	3.9 ± 0.0	2.1 ± 0.0	4.1 ± 0.1	2.5 ± 0.0	4.4 ± 0.0	2.0 ± 0.1	3.1 ± 0.1	2.1 ± 0.0	3.8 ± 0.1	2.2 ± 0.0
IPIP-NEO GPT-40-Min	ni											
Prompt-Demo	4.2 ± 0.0	2.9 ± 0.1	4.2 ± 0.1	3.2 ± 0.1	4.0 ± 0.0	2.6 ± 0.1	4.6 ± 0.1	2.4 ± 0.1	3.4 ± 0.0	2.1 ± 0.1	4.1 ± 0.0	2.6 ± 0.1
Prompt-Inst	4.8 ± 0.0	1.9 ± 0.2	4.9 ± 0.0	1.4 ± 0.0	4.9 ± 0.0	1.6 ± 0.0	4.8 ± 0.0	2.1 ± 0.1	4.9 ± 0.0	1.1 ± 0.1	4.9 ± 0.0	1.6 ± 0.1

Table 10: Full personality test results for various alignment methods, complementing Table 2.
Prompt-Demo-Sampling involves randomly sampling 10 examples from the entire BIG5-CHAT
dataset for each run, instead of using a fixed set of 10 random examples across runs. SFTPrompt-Demo and DPO-Prompt-Demo represent demonstration-based prompting applied to SFT
and DPO-trained models, respectively. Results for GPT-40-mini are presented in separate sections of the table. Scores range from 1 to 5, where a score closer to 5 indicates stronger agreement with the trait, while a score closer to 1 reflects weaker or opposing agreement.

1350 Figure 2 presents the BFI and IPIP-NEO test score results for the LLaMA-3 Instruct models, 1351 evaluated in zero-shot inference without any induced personality traits. The crowd-sourced re-1352 sponse scores for the BFI test are sourced from Huang et al. (2024), and those for the IPIP-1353 NEO test are drawn from Jiang et al. (2023). The results indicate that the scores for both 1354 LLaMA-3-8B-Instruct and LLaMA-3-70B-Instruct fall within the standard deviation of the human distribution. However, while LLaMA-3-8B-Instruct tends to generate more neutral 1355 scores (around 3 across most of the Big Five traits), LLaMA-3-70B-Instruct exhibits higher 1356 scores for openness, conscientiousness, extraversion, and agreeableness, and lower scores for neu-1357 roticism. 1358



Figure 2: The personality test results for the crowd and the LLaMA-3-Instruct models were obtained using zero-shot inference without explicitly inducing personality traits. The BFI test scores are displayed on the left. The IPIP-NEO test scores are displayed on the right.

1372 1373 1374

1375

1370

1371

C.4 INTRA-TRAIT CORRELATIONS IN PERSONALITY ASSESSMENT

1376 To assess how well the prompting and training methods simulate intra-trait correlations observed in human data, we first calculated the intra-trait correlations from real human distributions using 1377 the IPIP-NEO questionnaire, based on the PAPI-120-600K dataset from Zhu et al. (2024), which 1378 includes 619K human responses to the IPIP-NEO. Next, we computed the intra-trait correlations 1379 for the prompting, SFT, and DPO methods using the results from Table 2. These correlations are 1380 visualized in Figure 3, showing that most traits are positively correlated, with the exception of neu-1381 roticism. To quantify the similarity between the method-generated and human correlation matrices, 1382 we calculated the matrix distance using the Frobenius norm, where 0 represents perfect similarity 1383 and 10 indicates maximum dissimilarity. The matrix distances were 2.10 for prompting, 1.55 for 1384 SFT, and 2.06 for DPO. These results suggest that the trained models, particularly SFT, more ac-1385 curately capture the trait correlations seen in natural human data compared to the prompting-based 1386 methods.

1387 1388

1388 C.5 REASONING BENCHMARK RESULTS FOR LLAMA-3-70B-INSTRUCT 1389

The complete results for the general reasoning tasks evaluated on the LLaMA-3-70B-Instruct model are presented in Table 11. Note that the GPQA results in Table 3 were obtained using zero-shot prompting. This evaluation encompasses multiple reasoning domains and highlights the impact of different training methodologies: prompting, SFT, and DPO. These methods were assessed based on their ability to preserve the reasoning capabilities.

The results indicate that the SFT method consistently delivers the strongest performance across the benchmarks, outperforming both DPO and the prompting-based approach. For the 70B model, SFT emerges as the most effective method, achieving an optimal balance between incorporating personality traits and maintaining robust reasoning functionality. The aggregated results underscore the reliability of SFT, which demonstrates superior performance across diverse reasoning tasks, making it a robust choice for large-scale language models.

In contrast, the performance of the DPO method is more variable. While DPO excels in certain
 scenarios, such as low Neuroticism within the TruthfulQA task—where it achieves a notable score
 of 65.8%—its overall results are less consistent across other reasoning benchmarks. Moreover, the
 final average scores reveal that high-trait DPO models underperform compared to their low-trait



Figure 3: Intra-trait Pearson correlations for human distributions on IPIP-NEO and the correspond-1414 ing results from instruction-based prompting, SFT, and DPO. O represents openness, C conscien-1415 tiousness, E extraversion, A agreeableness, and N neuroticism. The correlations especially for SFT 1416 align well with human distributions across openness, conscientiousness, extraversion, and agreeable-1417 ness. Neuroticism shows less alignment with the other four traits compared to human distribution.

Benchmark	Direct	Method	Oper High	nness Low	Conscien High	tiousness Low	Extrav High	version Low	Agreea High	bleness Low	Neuro High	ticism Low	Ave High	rage Low
Hallucination Detec	tion													
TruthfulQA	58.6 ± 1.7	Prompt SFT DPO	$\begin{array}{c} 54.1 \pm 1.6 \\ 55.2 \pm 1.6 \\ 54.6 \pm 1.6 \end{array}$	$\begin{array}{c} 51.1 \pm 1.6 \\ 52.8 \pm 1.6 \\ 54.2 \pm 1.7 \end{array}$	$\begin{array}{c} 55.9 \pm 1.7 \\ 55.6 \pm 1.6 \\ 64.6 \pm 1.6 \end{array}$	$\begin{array}{c} 45.2 \pm 1.6 \\ 50.8 \pm 1.5 \\ 38.5 \pm 1.6 \end{array}$	$\begin{array}{c} 52.0 \pm 1.6 \\ 54.5 \pm 1.6 \\ 46.0 \pm 1.7 \end{array}$	$\begin{array}{c} 55.7 \pm 1.6 \\ 56.7 \pm 1.6 \\ 65.3 \pm 1.6 \end{array}$	$\begin{array}{c} 52.3 \pm 1.7 \\ 54.4 \pm 1.6 \\ 59.6 \pm 1.6 \end{array}$	$\begin{array}{c} 49.1 \pm 1.6 \\ 51.6 \pm 1.6 \\ 50.6 \pm 1.6 \end{array}$	$\begin{array}{c} 48.9 \pm 1.6 \\ 52.4 \pm 1.5 \\ 43.0 \pm 1.7 \end{array}$	$\begin{array}{c} 58.6 \pm 1.6 \\ 56.7 \pm 1.6 \\ 65.8 \pm 1.6 \end{array}$	$\begin{array}{c} 52.6 \pm 1.6 \\ 54.4 \pm 1.6 \\ 53.6 \pm 1.6 \end{array}$	$\begin{array}{c} 51.9 \pm 1.6 \\ 53.7 \pm 1.6 \\ 54.9 \pm 1.6 \end{array}$
Social Reasoning														
SocialIQA	46.6 ± 1.1	Prompt SFT DPO	$\begin{array}{c} 40.8\pm1.1\\ 50.3\pm1.1\\ 41.5\pm1.1 \end{array}$	$\begin{array}{c} 43.9 \pm 1.1 \\ 50.4 \pm 1.1 \\ 44.5 \pm 1.1 \end{array}$	$\begin{array}{c} 42.9 \pm 1.1 \\ 50.9 \pm 1.1 \\ 44.7 \pm 1.1 \end{array}$	$\begin{array}{c} 39.9 \pm 1.1 \\ 46.8 \pm 1.1 \\ 37.6 \pm 1.1 \end{array}$	$\begin{array}{c} 43.3 \pm 1.1 \\ 50.0 \pm 1.1 \\ 43.0 \pm 1.1 \end{array}$	$\begin{array}{c} 42.0\pm1.1\\ 50.3\pm1.1\\ 43.6\pm1.1\end{array}$	$\begin{array}{c} 42.4 \pm 1.1 \\ 50.5 \pm 1.1 \\ 44.8 \pm 1.1 \end{array}$	$\begin{array}{c} 40.8 \pm 1.1 \\ 46.6 \pm 1.1 \\ 39.0 \pm 1.1 \end{array}$	$\begin{array}{c} 39.1 \pm 1.1 \\ 48.2 \pm 1.1 \\ 40.0 \pm 1.1 \end{array}$	$\begin{array}{c} 44.1 \pm 1.1 \\ 50.6 \pm 1.1 \\ 45.3 \pm 1.1 \end{array}$	$\begin{array}{c} 41.7\pm1.1\\ 50.0\pm1.1\\ 42.8\pm1.1 \end{array}$	$\begin{array}{c} 42.1 \pm 1.1 \\ 48.9 \pm 1.1 \\ 42.0 \pm 1.1 \end{array}$
Commonsense Reas	oning													
CommonsenseQA	27.0 ± 1.3	Prompt SFT DPO	$\begin{array}{c} 60.0 \pm 1.4 \\ 77.7 \pm 1.2 \\ 57.7 \pm 1.4 \end{array}$	$\begin{array}{c} 59.9 \pm 1.4 \\ 78.8 \pm 1.2 \\ 65.9 \pm 1.4 \end{array}$	$\begin{array}{c} 22.5 \pm 1.2 \\ 77.6 \pm 1.2 \\ 23.8 \pm 1.2 \end{array}$	$\begin{array}{c} 22.3 \pm 1.2 \\ 66.0 \pm 1.4 \\ 25.8 \pm 1.3 \end{array}$	$\begin{array}{c} 35.5 \pm 1.4 \\ 75.7 \pm 1.2 \\ 23.2 \pm 1.2 \end{array}$	$\begin{array}{c} 50.0 \pm 1.4 \\ 78.9 \pm 1.2 \\ 70.8 \pm 1.3 \end{array}$	$\begin{array}{c} 45.0 \pm 1.4 \\ 77.0 \pm 1.2 \\ 21.3 \pm 1.2 \end{array}$	$\begin{array}{c} 34.9 \pm 1.4 \\ 73.8 \pm 1.3 \\ 39.2 \pm 1.4 \end{array}$	$\begin{array}{c} 20.2 \pm 1.2 \\ 79.1 \pm 1.2 \\ 20.1 \pm 1.1 \end{array}$	$\begin{array}{c} 36.8 \pm 1.4 \\ 78.5 \pm 1.2 \\ 44.6 \pm 1.4 \end{array}$	$\begin{array}{c} 36.6 \pm 1.3 \\ 77.4 \pm 1.2 \\ 29.2 \pm 1.2 \end{array}$	$\begin{array}{c} 40.8 \pm 1.4 \\ 75.2 \pm 1.3 \\ 49.3 \pm 1.4 \end{array}$
PIQA	80.4 ± 0.9	Prompt SFT DPO	$\begin{array}{c} 79.6 \pm 0.9 \\ 81.2 \pm 0.9 \\ 76.4 \pm 1.0 \end{array}$	$\begin{array}{c} 79.8 \pm 0.9 \\ 81.0 \pm 0.9 \\ 76.8 \pm 1.0 \end{array}$	$\begin{array}{c} 80.5\pm 0.9\\ 81.2\pm 0.9\\ 79.4\pm 0.9\end{array}$	$\begin{array}{c} 77.3 \pm 1.0 \\ 80.4 \pm 0.9 \\ 70.9 \pm 1.1 \end{array}$	$\begin{array}{c} 78.0 \pm 1.0 \\ 81.8 \pm 0.9 \\ 76.4 \pm 1.0 \end{array}$	$\begin{array}{c} 80.0\pm 0.9\\ 81.3\pm 0.9\\ 79.8\pm 0.9\end{array}$	$\begin{array}{c} 79.8 \pm 0.9 \\ 81.2 \pm 0.9 \\ 78.5 \pm 1.0 \end{array}$	$\begin{array}{c} 78.4 \pm 1.0 \\ 80.0 \pm 0.9 \\ 74.0 \pm 1.0 \end{array}$	$\begin{array}{c} 78.8 \pm 1.0 \\ 81.0 \pm 0.9 \\ 72.9 \pm 1.0 \end{array}$	$\begin{array}{c} 80.7\pm 0.9\\ 81.2\pm 0.9\\ 79.5\pm 0.9\end{array}$	$\begin{array}{c} 79.3 \pm 0.9 \\ 81.3 \pm 0.9 \\ 76.7 \pm 1.0 \end{array}$	$\begin{array}{c} 79.2 \pm 0.9 \\ 80.8 \pm 0.9 \\ 76.2 \pm 1.0 \end{array}$
Math Reasoning														
GSM8K	80.6 ± 1.1	Prompt SFT DPO	$\begin{array}{c} 75.7 \pm 1.2 \\ 85.8 \pm 1.0 \\ 87.9 \pm 0.9 \end{array}$	$\begin{array}{c} 70.1 \pm 1.3 \\ 76.2 \pm 1.2 \\ 88.5 \pm 0.9 \end{array}$	$\begin{array}{c} 73.5 \pm 1.2 \\ 86.4 \pm 0.9 \\ 90.2 \pm 0.8 \end{array}$	$\begin{array}{c} 32.6 \pm 1.3 \\ 81.7 \pm 1.1 \\ 80.6 \pm 1.1 \end{array}$	$\begin{array}{c} 80.8 \pm 1.1 \\ 85.1 \pm 1.0 \\ 88.9 \pm 0.9 \end{array}$	$\begin{array}{c} 33.5 \pm 1.3 \\ 86.7 \pm 0.9 \\ 90.4 \pm 0.8 \end{array}$	$\begin{array}{c} 87.2 \pm 0.9 \\ 87.0 \pm 0.9 \\ 87.3 \pm 0.9 \end{array}$	$\begin{array}{c} 77.8 \pm 1.1 \\ 74.5 \pm 1.2 \\ 90.0 \pm 0.8 \end{array}$	$\begin{array}{c} 26.0 \pm 1.2 \\ 76.0 \pm 1.2 \\ 15.2 \pm 1.0 \end{array}$	$\begin{array}{c} 89.4 \pm 0.8 \\ 87.3 \pm 0.9 \\ 91.0 \pm 0.8 \end{array}$	$\begin{array}{c} 68.6 \pm 1.1 \\ 84.1 \pm 1.0 \\ 73.9 \pm 0.9 \end{array}$	$\begin{array}{c} 60.7 \pm 1.2 \\ 81.3 \pm 1.1 \\ 88.1 \pm 0.9 \end{array}$
MathQA	39.0 ± 0.9	Prompt SFT DPO	$\begin{array}{c} 33.5\pm0.9\\ 43.3\pm0.9\\ 33.9\pm0.9 \end{array}$	$\begin{array}{c} 33.5\pm 0.9\\ 42.6\pm 0.9\\ 34.7\pm 0.9\end{array}$	$\begin{array}{c} 32.8 \pm 0.9 \\ 43.0 \pm 0.9 \\ 32.9 \pm 0.9 \end{array}$	$\begin{array}{c} 31.5\pm0.9\\ 43.3\pm0.9\\ 28.1\pm0.8 \end{array}$	$\begin{array}{c} 32.3 \pm 0.9 \\ 43.2 \pm 0.9 \\ 30.5 \pm 0.8 \end{array}$	$\begin{array}{c} 33.3 \pm 0.9 \\ 42.7 \pm 0.9 \\ 35.0 \pm 0.9 \end{array}$	$\begin{array}{c} 33.6 \pm 0.9 \\ 42.9 \pm 0.9 \\ 31.3 \pm 0.8 \end{array}$	$\begin{array}{c} 32.4 \pm 0.9 \\ 42.9 \pm 0.9 \\ 32.8 \pm 0.9 \end{array}$	$\begin{array}{c} 32.1 \pm 0.9 \\ 42.8 \pm 0.9 \\ 28.9 \pm 0.8 \end{array}$	$\begin{array}{c} 34.1 \pm 0.9 \\ 43.3 \pm 0.9 \\ 34.0 \pm 0.9 \end{array}$	$\begin{array}{c} 32.9\pm 0.9\\ 43.0\pm 0.9\\ 31.5\pm 0.8 \end{array}$	$\begin{array}{c} 33.0\pm 0.9\\ 43.0\pm 0.9\\ 32.9\pm 0.9\end{array}$
General Reasoning														
MMLU	74.5 ± 0.3	Prompt SFT DPO	$\begin{array}{c} 70.3 \pm 0.4 \\ 72.5 \pm 0.4 \\ 57.9 \pm 0.4 \end{array}$	$\begin{array}{c} 69.6 \pm 0.4 \\ 72.0 \pm 0.4 \\ 64.4 \pm 0.4 \end{array}$	$\begin{array}{c} 40.6 \pm 0.4 \\ 73.1 \pm 0.4 \\ 50.3 \pm 0.4 \end{array}$	$\begin{array}{c} 52.8 \pm 0.4 \\ 68.6 \pm 0.4 \\ 33.8 \pm 0.4 \end{array}$	$\begin{array}{c} 56.9 \pm 0.4 \\ 72.1 \pm 0.4 \\ 42.3 \pm 0.4 \end{array}$	$\begin{array}{c} 72.8 \pm 0.4 \\ 73.5 \pm 0.4 \\ 72.3 \pm 0.4 \end{array}$	$\begin{array}{c} 69.0 \pm 0.4 \\ 72.8 \pm 0.4 \\ 34.3 \pm 0.4 \end{array}$	$\begin{array}{c} 69.2 \pm 0.4 \\ 70.7 \pm 0.4 \\ 62.5 \pm 0.4 \end{array}$	$\begin{array}{c} 55.3 \pm 0.4 \\ 72.5 \pm 0.4 \\ 33.2 \pm 0.4 \end{array}$	$\begin{array}{c} 67.9 \pm 0.4 \\ 73.8 \pm 0.4 \\ 69.1 \pm 0.4 \end{array}$	$\begin{array}{c} 58.4 \pm 0.4 \\ 72.6 \pm 0.4 \\ 43.6 \pm 0.4 \end{array}$	$\begin{array}{c} 66.5 \pm 0.4 \\ 71.7 \pm 0.4 \\ 60.4 \pm 0.4 \end{array}$
GPQA (0-shot)	33.5 ± 2.2	Prompt SFT DPO	$\begin{array}{c} 31.5 \pm 2.2 \\ 33.5 \pm 2.2 \\ 36.8 \pm 2.3 \end{array}$	$\begin{array}{c} 34.2 \pm 2.2 \\ 32.4 \pm 2.2 \\ 31.9 \pm 2.2 \end{array}$	$\begin{array}{c} 31.7 \pm 2.2 \\ 34.2 \pm 2.2 \\ 35.7 \pm 2.3 \end{array}$	$\begin{array}{c} 32.4 \pm 2.2 \\ 34.2 \pm 2.2 \\ 30.6 \pm 2.2 \end{array}$	$\begin{array}{c} 34.6 \pm 2.2 \\ 33.3 \pm 2.2 \\ 35.9 \pm 2.3 \end{array}$	$\begin{array}{c} 32.1 \pm 2.2 \\ 34.4 \pm 2.2 \\ 35.9 \pm 2.3 \end{array}$	$\begin{array}{c} 32.4 \pm 2.2 \\ 33.3 \pm 2.2 \\ 35.5 \pm 2.3 \end{array}$	$\begin{array}{c} 32.8 \pm 2.2 \\ 33.3 \pm 2.2 \\ 35.7 \pm 2.3 \end{array}$	$\begin{array}{c} 31.9 \pm 2.2 \\ 34.4 \pm 2.2 \\ 32.6 \pm 2.2 \end{array}$	$\begin{array}{c} 32.1 \pm 2.2 \\ 33.5 \pm 2.2 \\ 34.6 \pm 2.2 \end{array}$	$\begin{array}{c} 32.4 \pm 2.2 \\ 33.7 \pm 2.2 \\ 35.3 \pm 2.3 \end{array}$	$\begin{array}{c} 32.7 \pm 2.2 \\ 33.6 \pm 2.2 \\ 33.7 \pm 2.2 \end{array}$
GPQA (5-shot)	36.6 ± 2.3	Prompt SFT DPO	$\begin{array}{c} 35.9 \pm 2.3 \\ 32.4 \pm 2.2 \\ 37.5 \pm 2.3 \end{array}$	$\begin{array}{c} 32.6 \pm 2.2 \\ 32.8 \pm 2.2 \\ 31.2 \pm 2.2 \end{array}$	$\begin{array}{c} 36.2 \pm 2.3 \\ 34.4 \pm 2.2 \\ 35.9 \pm 2.3 \end{array}$	$\begin{array}{c} 35.7 \pm 2.3 \\ 33.7 \pm 2.2 \\ 31.2 \pm 2.2 \end{array}$	$\begin{array}{c} 36.2 \pm 2.3 \\ 33.0 \pm 2.2 \\ 37.1 \pm 2.3 \end{array}$	$\begin{array}{c} 35.7 \pm 2.3 \\ 33.9 \pm 2.2 \\ 35.5 \pm 2.3 \end{array}$	$\begin{array}{c} 34.4 \pm 2.2 \\ 33.7 \pm 2.2 \\ 33.5 \pm 2.2 \end{array}$	$\begin{array}{c} 34.8 \pm 2.3 \\ 32.8 \pm 2.2 \\ 32.1 \pm 2.2 \end{array}$	$\begin{array}{c} 36.6 \pm 2.3 \\ 33.7 \pm 2.2 \\ 36.6 \pm 2.3 \end{array}$	$\begin{array}{c} 34.2\pm2.2\\ 34.8\pm2.3\\ 35.7\pm2.3 \end{array}$	$\begin{array}{c} 35.9 \pm 2.3 \\ 33.4 \pm 2.2 \\ 36.1 \pm 2.3 \end{array}$	$\begin{array}{c} 34.6 \pm 2.3 \\ 33.6 \pm 2.2 \\ 33.1 \pm 2.2 \end{array}$
Average	53.0 ± 1.3	Prompt SFT DPO	$\begin{array}{c} 53.5 \pm 1.3 \\ 59.1 \pm 1.3 \\ 53.8 \pm 1.3 \end{array}$	$\begin{array}{c} 52.7 \pm 1.3 \\ 57.7 \pm 1.3 \\ 54.7 \pm 1.3 \end{array}$	$\begin{array}{c} 46.3 \pm 1.3 \\ 59.6 \pm 1.3 \\ 50.8 \pm 1.3 \end{array}$	$\begin{array}{c} 41.1 \pm 1.3 \\ 56.2 \pm 1.3 \\ 41.9 \pm 1.3 \end{array}$	$\begin{array}{c} 50.0 \pm 1.3 \\ 58.7 \pm 1.3 \\ 47.0 \pm 1.3 \end{array}$	$\begin{array}{r} 48.3 \pm 1.3 \\ 59.8 \pm 1.3 \\ 58.7 \pm 1.3 \end{array}$	$\begin{array}{c} 52.9 \pm 1.3 \\ 59.2 \pm 1.3 \\ 47.3 \pm 1.3 \end{array}$	$\begin{array}{c} 50.0 \pm 1.3 \\ 56.2 \pm 1.3 \\ 50.7 \pm 1.3 \end{array}$	$\begin{array}{c} 41.0 \pm 1.3 \\ 57.8 \pm 1.3 \\ 35.8 \pm 1.3 \end{array}$	$\begin{array}{c} 53.1 \pm 1.3 \\ 60.0 \pm 1.3 \\ 55.5 \pm 1.3 \end{array}$	$\begin{array}{c} 48.7 \pm 1.3 \\ 58.9 \pm 1.3 \\ 47.0 \pm 1.3 \end{array}$	$\begin{array}{c} 49.1 \pm 1.3 \\ 58.0 \pm 1.3 \\ 52.3 \pm 1.3 \end{array}$

Table 11: Benchmark results for different personality traits on LLaMA-3-70B-Instruct. Di-1439 rect refers to direct inference without including personality-related prompts. Prompt refers to 1440 instruction-based prompting. The table includes standard errors (shown as \pm values) to provide 1441 statistical context for the results. 1442

1444 counterparts in general. This suggests a potential misalignment between DPO's training objectives 1445 and the reasoning requirements of specific tasks. These findings highlight the nuanced trade-offs 1446 between training strategies, with SFT offering the most reliable approach for balancing personality 1447 trait integration and cognitive task performance in large-scale models.

1448 1449 1450

1438

1443

1418

C.6 **REASONING BENCHMARK RESULTS FOR** LLAMA-3-8B-INSTRUCT

The reasoning evaluation results for the LLaMA-3-8B-Instruct model, assessed across six rea-1451 soning domains, are summarized in Table 12. Overall, the DPO method generally outperformed SFT 1452 and demonstrated performance comparable to the prompt-based approach. This indicates that, with 1453 the smaller 8B model, DPO effectively aligns personality traits without significantly compromising 1454 reasoning capabilities. 1455

A comparison of personality trait levels revealed that models simulating high trait levels consistently 1456 outperformed their low-trait counterparts in both DPO and SFT settings. For instance, on the Truth-1457 fulQA benchmark, the high-conscientiousness DPO model achieved 55.0%, significantly surpassing

the low-conscientiousness model's 39.0%. Similarly, on the GSM8K math reasoning task, the high-conscientiousness DPO model scored 72.2%, substantially outperforming the low-level model.

On benchmarks such as TruthfulQA, GPQA (both zero-shot and five-shot), and MathQA, models trained using SFT and DPO performed comparably to the original unaligned model. This suggests that personality trait alignment does not adversely affect reasoning performance in these tasks for a small model. However, notable variations were observed in other benchmarks. For example, DPO exhibited significantly reduced performance on CommonsenseQA and MMLU compared to SFT, prompting, and the original model. Conversely, SFT underperformed on the GSM8K benchmark relative to DPO, prompting, and the original model. These results suggest that the DPO method may be more effective than SFT in preserving or enhancing reasoning performance for specific tasks and traits on small models, though the choice of alignment method may depend on the specific reasoning domain.

Benchmark	Original	Method	Oper High	nness Low	Consci High	entiousness Low	Extra High	version Low	Agree: High	ableness Low	Neuro High	ticism Low	Ave High	rage Low
Hallucination Dete	ction													
TruthfulQA	53.5	Prompt SFT	49.0 50.0	51.5 45.7	50.6 50.9	44.4 43.8 20.0	45.3 46.2	51.9 52.0	49.2 49.9	50.3 46.3	54.6 53.6	45.2 42.9	49.7 50.1	48.7 46.1
Code Reasoning		DFO	32.4	49.1	33.0	39.0	33.0	39.2	32.0	43.5	38.2	30.0	50.7	40.5
HumanEval	60.4	Prompt	59.1	59.8	62.2	61.6	61.0	63.4	62.8	62.2	60.4	61.6	61.1	61.7
		SFT DPO	57.9 57.3	54.3 0.6	59.8 27.4	56.1 0.0	58.5 43.3	57.3 0.0	60.4 8.5	54.9 32.9	58.5 0.0	58.5 7.9	59.0 27.3	56.2 8.3
MBPP	54.6	Prompt SFT DPO	54.6 56.2 53.6	55.4 56.2 47.6	54.2 54.2 53.0	55.2 56.2 35.2	55.8 56.4 54.6	56.0 56.4 51.4	55.4 55.6 54.4	54.8 55.8 53.8	54.4 55.0 52.0	55.8 56.4 54.2	54.9 55.5 42.9	55.4 56.2 48.4
Social Reasoning		DIG	55.0	47.0	55.0	55.2	54.0	51.4	51.1	55.0	52.0	51.2	12.9	-10.1
SocialIQA	49.7	Prompt SFT DPO	41.9 44.0 43.8	42.3 44.9 43.8	41.1 45.9 42.5	39.3 41.9 37.8	41.5 44.4 41.8	41.6 44.6 40.9	41.8 43.7 42.8	39.5 41.4 38.4	42.1 44.6 42.8	39.4 40.8 39.0	41.7 44.5 42.7	40.4 42.7 40.0
Commonsense Rea	soning													
CommonsenseQA	51.8	Prompt SFT DPO	64.6 61.8 22.9	60.6 57.9 24.8	38.0 50.5 48.2	31.3 34.3 21.6	45.9 52.7 29.1	55.0 60.8 56.6	55.4 55.4 28.4	36.3 36.0 26.3	33.9 63.4 47.7	23.3 30.6 23.7	47.6 56.8 35.3	41.3 43.9 30.6
Math Reasoning														
GSM8K	64.7	Prompt SFT DPO	13.5 19.8 68.4	58.4 0.5 31.8	23.4 20.2 72.2	61.0 1.4 31.8	40.0 6.0 69.7	57.1 0.5 63.0	29.3 6.4 70.7	71.6 4.8 64.8	24.1 20.1 71.9	31.9 53.3 3.0	26.1 14.5 70.6	56.0 12.1 38.9
MathQA	27.9	Prompt SFT DPO	27.6 30.1 26.9	28.3 30.2 27.8	27.9 29.6 28.3	27.3 30.3 25.1	27.1 31.0 25.8	27.8 30.6 27.6	27.2 29.6 24.9	28.1 30.3 27.7	28.1 29.6 29.7	25.9 29.4 24.9	27.6 30.0 27.1	27.5 30.2 26.6
General Knowledge	?													
MMLU	51.2	Prompt SFT DPO	37.5 45.0 23.0	29.1 48.5 29.8	23.2 35.6 29.7	27.0 32.0 26.9	24.7 37.5 24.8	29.2 46.5 41.4	27.7 44.2 30.7	25.5 39.9 26.3	23.4 47.1 30.8	23.8 31.7 23.1	27.3 41.9 27.8	26.9 39.7 29.5
GPQA (0-shot)	28.1	Prompt SFT DPO	29.0 27.9 27.9	28.8 27.9 25.0	28.6 28.1 29.7	23.0 25.0 21.0	28.6 27.2 27.2	29.2 28.3 26.8	29.0 28.8 28.8	27.2 24.1 21.4	28.8 29.0 29.5	28.3 28.3 25.2	28.8 28.2 28.6	27.3 26.7 23.9
GPQA (5-shot)	29.9	Prompt SFT DPO	29.7 26.1 27.9	26.6 27.0 26.3	28.8 28.8 28.3	26.8 26.6 23.0	28.3 28.8 26.8	26.6 28.6 28.1	27.9 30.6 27.5	28.6 27.9 24.6	29.0 28.6 28.8	25.2 27.5 25.2	28.7 28.6 27.9	26.8 27.5 25.4
Average	43.9	Prompt SFT DPO	35.8 37.2 35.6	40.5 34.0 30.7	31.5 34.8 41.6	34.4 27.6 26.9	34.3 32.8 34.1	39.5 35.3 43.2	35.1 35.0 37.7	38.2 29.9 33.8	31.7 38.8 42.4	29.1 34.8 23.4	33.7 35.7 38.3	36.4 32.3 31.6

Table 12: Benchmark results for the LLaMA-3-8B-Instruct model are presented across vari-ous personality traits and evaluation methods. The benchmarks are categorized into six key areas: Hallucination Detection, General Reasoning, Social Reasoning, Commonsense Reasoning, Mathe-matical Reasoning, and General Knowledge.

1512 D CORRELATION BETWEEN PERSONALITY TRAITS AND REASONING BEHAVIORS

1515

1517

1531

1516 D.1 HUMAN VS. LLAMA-3-70B-INSTRUCT

Understanding the influence of personality traits on reasoning behaviors in LLMs is crucial for developing models tailored to specific personality profiles. Research on the Big Five personality traits has consistently demonstrated their significant impact on human cognition and problem-solving abilities (John et al., 1999; Soto et al., 2011). Traits such as openness, conscientiousness, and agreeableness are often associated with enhanced reasoning capabilities, while neuroticism has been found to impair performance across a range of reasoning tasks (Ackerman & Heggestad, 1997; Schaie et al., 2004; Chamorro-Premuzic et al., 2006).

1525Table 13 summarizes relevant findings from recent psychological studies and their alignment with
our experimental results on LLaMA-3-70B-Instruct. Our findings corroborate these stud-
ies, indicating that models exhibiting higher conscientiousness and agreeableness generally perform
better in reasoning tasks. In contrast, models characterized by lower levels of extraversion and neu-
roticism also demonstrate improved reasoning performance. These results highlight the potential of
personality-aligned training to optimize LLM performance for reasoning-intensive tasks.

1532 **Openness** Openness is associated with intellectual curiosity and creativity and enhances 1533 problem-solving in tasks requiring abstract reasoning and social cognition (Ackerman & 1534 Heggestad, 1997; McCrae, 1987). While research indicates that openness positively corre-1535 lates with cognitive abilities (Chamorro-Premuzic et al., 2006; Costa Jr et al., 1976; Graham 1536 & Lachman, 2012; Schaie et al., 2004), our models do not show significant performance dif-1537 ferences across reasoning tasks based on openness levels, with the exception of SFT on math 1538 reasoning tasks. This suggests that openness may not directly translate to gains in reasoning 1539 tasks beyond math, despite its known benefits to human cognition.

Conscientiousness Conscientiousness, linked to discipline and organization, consistently improves model performance in mathematical reasoning and hallucination detection. This aligns with psychological studies showing that higher conscientiousness is linked to better academic performance and fewer errors in cognitive tasks due to increased diligence and thoroughness (Roberts et al., 2014; Poropat, 2009; Digman, 1990; Moutafi et al., 2003; Schaie et al., 2004).

Extraversion Extraversion is often associated with sociability and shows mixed results in cognitive tasks. While it can enhance social reasoning, it may negatively affect individual problem-solving tasks, such as math reasoning (Blickle, 1996; Ashton et al., 2002; Costa Jr et al., 1976). Our models simulating lower extraversion perform better across many reasoning domains, including math and also commonsense reasoning, consistent with findings that high extraversion can detract from tasks requiring focused, solitary work (Matthews & Gilliland, 1999; Chamorro-Premuzic & Furnham, 2006).

Agreeableness Agreeableness, linked to traits like trust and cooperation, improves social reasoning in our models, consistent with human studies (Graziano, 1997). However, it shows minimal impact on math or commonsense reasoning, reflecting research suggesting that agreeableness is less beneficial for analytical tasks (Poropat, 2009; Ackerman & Heggestad, 1997; Schaie et al., 2004).

Neuroticism Neuroticism reflects emotional instability, and is consistently associated with poorer cognitive performance due to anxiety and cognitive interference, especially social reasoning and hallucination detection (Robinson & Tamir, 2005; Zeidner, 2005; Chamorro-Premuzic et al., 2006; Eysenck, 2013). Our models confirm this, with lower Neuroticism levels leading to better performance across almost all reasoning tasks.

1562

1552

1553

1554

1555

1556

Table 13: Summary of the influence of Big Five personality traits on reasoning tasks in human cognition, and comparison of psychological research findings with our experimental results on LLMs.

1566 D.2 HUMAN VS. LLAMA-3-8B-INSTRUCT

The influence of Big Five Personality traits on reasoning tasks in human cognition, as outlined in Table 13, served as a foundation for analyzing the performance of the LLaMA-3-8B-Instruct model. This analysis aims to explore how alignment with different personality traits affects the model's reasoning capabilities. Below, we summarize the observed correlations between each trait and the model's performance across various reasoning benchmarks.

Openness The impact of Openness on reasoning performance was highly task-dependent. Mod els aligned with high levels of Openness using the DPO method exhibited significantly improved
 performance in mathematical reasoning tasks. However, these models underperformed in common sense reasoning benchmarks compared to both the prompt-based approach and the original model.
 These results suggest that while high Openness alignment enhances mathematical reasoning, it does not guarantee consistent improvements across all reasoning domains.

1579 Conscientiousness A strong positive correlation was observed between Conscientiousness and rea 1580 soning performance. Models aligned with higher levels of Conscientiousness consistently outper 1581 formed their low-level counterparts across most benchmarks. This trend highlights that high Con 1582 scientiousness alignment likely enhances systematic reasoning and attention to detail, benefiting
 1583 performance across diverse reasoning tasks.

Extraversion Lower levels of Extraversion were associated with better performance across reasoning tasks. Specifically, in commonsense reasoning benchmarks, models with low Extraversion significantly outperformed those with high Extraversion. This negative correlation suggests that high Extraversion may introduce distractibility, potentially impeding performance in tasks that require focused attention and analytical reasoning.

 Agreeableness The influence of Agreeableness on reasoning performance was minimal and inconsistent. No clear advantage was observed for models aligned with either high or low levels of Agreeableness across the benchmarks. These findings indicate that Agreeableness has a weak correlation with the model's reasoning capabilities, suggesting its alignment has little effect on overall performance.

Neuroticism The relationship between Neuroticism and reasoning performance was inconsistent and did not align with expectations from human cognition studies. High Neuroticism models performed well in some reasoning tasks, while low Neuroticism models scored poorly in others. These results imply that high Neuroticism alignment does not necessarily impair reasoning performance, contrasting with psychological findings in humans. This discrepancy may arise from limitations in how Neuroticism is modeled and represented in the training process.