
LLMs and Personalities: Inconsistencies Across Scales

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This study investigates the application of human psychometric assessments to
2 large language models (LLMs) to examine their consistency and malleability in
3 exhibiting personality traits. We administered the Big Five Inventory (BFI) and
4 the Eysenck Personality Questionnaire-Revised (EPQ-R) to various LLMs across
5 different model sizes and persona prompts. Our results reveal substantial variability
6 in responses due to question order shuffling, challenging the notion of a stable
7 LLM "personality." Larger models demonstrated more consistent responses, while
8 persona prompts significantly influenced trait scores. Notably, the assistant per-
9 sona led to more predictable scaling, with larger models exhibiting more socially
10 desirable and less variable traits. In contrast, non-conventional personas displayed
11 unpredictable behaviors, sometimes extending personality trait scores beyond the
12 typical human range. These findings have important implications for understand-
13 ing LLM behavior under different conditions and reflect on the consequences of
14 scaling.

15 1 Introduction

16 Large language models (LLMs) have demonstrated remarkable capabilities in natural language
17 processing tasks, often exhibiting human-like responses in various contexts [Brown et al., 2020]. As
18 these models become more sophisticated, questions arise about the extent to which they can emulate
19 human-like personality traits and the consistency of such behaviors. Understanding these aspects is
20 crucial for both the development of more effective AI systems and for addressing ethical concerns
21 surrounding their deployment.

22 Personality testing, a central element of psychological assessment in humans, offers a structured
23 approach to probing these questions in LLMs. By applying established psychometric instruments to
24 AI models, we can gain insights into their ability to consistently exhibit personality traits and how
25 these traits may be influenced by different prompting strategies and model architectures.

26 Recent work has begun to explore this area [Huang, 2024, La Cava et al., 2024], delineating the most
27 prevalent psychological traits in LLMs. However, Gupta et al. [2024] have raised important concerns
28 about the reliability of using self-assessment personality tests with LLMs, showing high sensitivity to
29 prompt wording and option ordering.

30 Building upon these efforts, our study specifically examines the consistency of personality traits
31 across different model sizes, and the malleability of these traits by persona prompts. The reliability
32 of these traits was studied across multiple runs, each with shuffled question orders.

33 2 Methods

34 We employed two widely-used personality assessments: the Big Five Inventory (BFI), which assesses
35 five broad personality traits (Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism)

36 [John and Srivastava, 1999], and the Eysenck Personality Questionnaire-Revised (EPQ-R), which
37 measures three personality dimensions (Psychoticism, Extraversion, Neuroticism) and includes a Lie
38 scale [Eysenck et al., 1985].

39 We tested multiple versions of two LLM families: LLaMA 3.1 (8b, 70b, and 405b parameter versions,
40 [Dubey et al., 2024]) and Gemma 2 (9b and 27b parameter versions, [Team et al., 2024]). To evaluate
41 the impact of persona prompting, we tested four different personas: an assistant (helpful AI), a
42 Buddhist monk, an individual with psychopathic traits, and an individual with schizophrenia.

43 Our testing procedure involved administering both BFI (44 questions) and EPQ-R (100 questions)
44 to each model and persona combination. Questions were asked in batches of 8 for BFI and 10 for
45 EPQ-R. To assess consistency, we shuffled the questions randomly for each run. We conducted
46 100 runs for each model-persona combination to generate distributions of scores. Responses were
47 collected as numerical scores (1-5 for BFI, 0 or 1 for EPQ-R), and we accounted for reverse-scored
48 items. We also included a baseline "random" condition where responses were generated randomly to
49 serve as a point of comparison. A detailed prompt contained the instructions on how to perform the
50 questionnaire (see Appendix).

51 Each persona was implemented using a specific prompt (preceding the instructions) describing the
52 characteristics and background of the persona (see Appendix). It is important to note that all LLMs,
53 including those in the "assistant" condition, were asked to take up a persona, as the concept of an AI
54 assistant itself represents a form of persona.

55 Part of the code used in this study was adapted from Huang [2024], with fixes and substantial
56 expansions made to suit the specific needs of our research design.

57 **3 Results**

58 Figure 1 presents the distribution of BFI scores across different models, personas, and traits. We
59 observed substantial variability in responses due to question order shuffling, particularly in smaller
60 models. Larger models showed more stable BFI trait scores across runs, with narrower distributions
61 compared to smaller models. This trend was particularly evident for the assistant persona. The impact
62 of different personas on personality profiles was significant and aligned with expected characteristics.
63 The assistant persona consistently scored high on Agreeableness and Conscientiousness, with low
64 variability across runs. The Buddhist monk persona exhibited high Openness and Agreeableness,
65 with remarkably low Neuroticism. The psychopathic traits persona showed low Agreeableness and
66 high Extraversion. The schizophrenia persona demonstrated high Neuroticism and low Extraversion.
67 Notably, the assistant persona led to more predictable scaling, with larger models exhibiting more
68 socially desirable (higher Agreeableness and Conscientiousness) and less variable traits. In contrast,
69 non-conventional personas displayed more unpredictable behaviors, sometimes extending personality
70 trait scores beyond the typical human range.

71 Figure 2 illustrates the distribution of EPQ-R scores for each trait across different models and
72 personas, corroborating and extending the findings from the BFI assessment. We observed substantial
73 variability in responses due to question order shuffling, particularly in smaller models. Notably,
74 larger models demonstrated more consistent responses across runs, as evidenced by tighter score
75 distributions. As observed in the BFI results, the trait scores for different personas were consistent
76 with the instructions given, reflecting the expected characteristics of each persona.

77 In both assessments, persona instructions contributed to reduced variability in certain cases, espe-
78 cially when a specific trait was clearly delineated in a specific profile. An observation across both
79 assessments was the strong tendency towards socially desirable responses in the assistant persona,
80 as evidenced by high Agreeableness and Conscientiousness in BFI, and low Psychoticism and high
81 Lying scores in EPQ-R.

82 **4 Discussion**

83 Our findings raise important questions about the nature of "personality" in LLMs and the interpretation
84 of their responses to psychological assessments. The high variability observed, especially in smaller
85 models, challenges the notion of a stable LLM personality and highlights the sensitivity of these

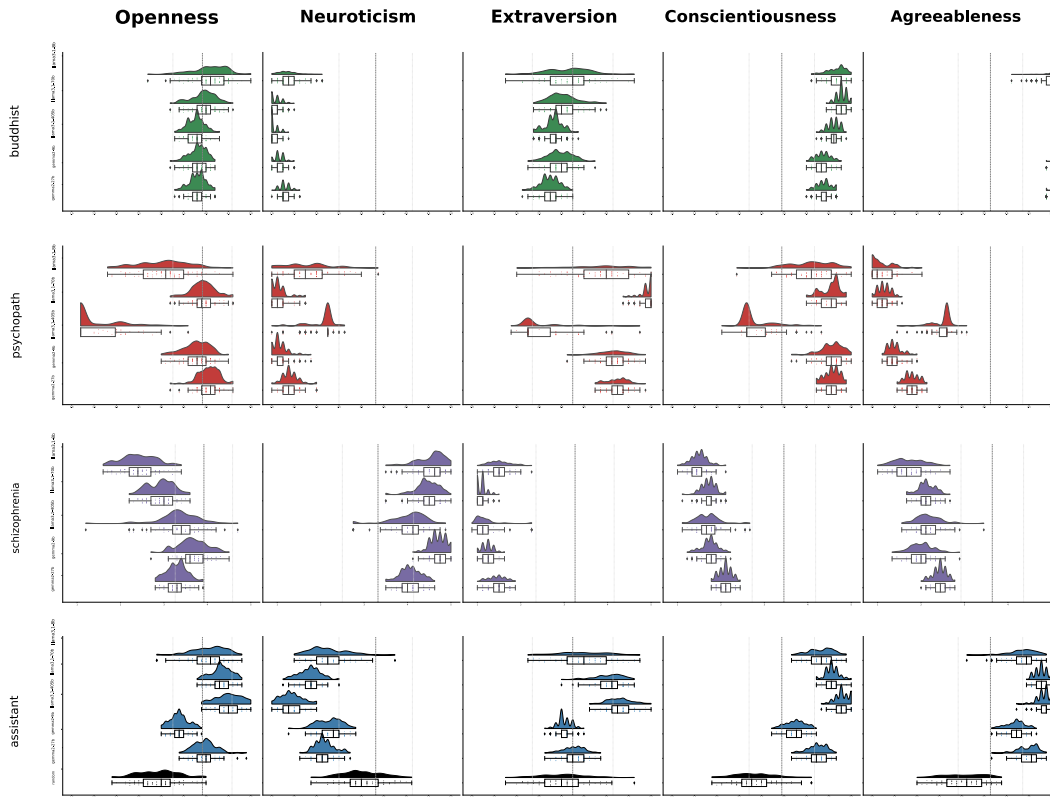


Figure 1: Distribution of BFI scores across different models, personas, and traits. Each violin plot represents the distribution of scores from 100 runs with randomly shuffled question order. The width of each violin indicates the density of scores at that level, with density values normalized within each quadrant. Inside each violin, a box plot shows the median (white dot), interquartile range (thick bar), and whiskers (thin lines). Individual points represent outliers. Colors represent different personas: Green - Buddhist, Red - Psychopath, Purple - Schizophrenia, Blue - Assistant. The solid black vertical line represents mean values for the human population, while the dashed line indicates the standard deviation of that mean. Models on the y-axis of each quadrant (from top to bottom) are: LLaMA 3.1 8b, 70b, 405b, and Gemma 2 9b, 27b. The bottom plot of the assistant quadrant shows a baseline condition labeled "random," representing scores generated by uniformly sampling responses (chance-level performance).

86 systems to input ordering. This variability suggests that caution should be exercised when attributing
 87 human-like personality traits to AI systems based on single interactions or assessments.

88 The relationship between model size and response consistency in the assistant persona suggests that
 89 larger models may develop more stable internal representations. This finding indicates that increased
 90 model capacity is necessary for more reliable and consistent helpful assistant personality emulation.
 91 However, it's crucial to note that even the largest models still exhibited variability.

92 Conversely, in the case of non-assistant personas, we observed U-shaped behaviors. This highlights
 93 an important consideration regarding the optimization of LLMs. While increasing model size and
 94 optimizing for benchmark performance may lead to monotonic increases in accuracy, our findings
 95 suggest that this may cause nonlinear shifts in personality traits for non-assistant personas. This
 96 observation could have implications for the deployment of AI systems requiring specific role-playing.

97 The effectiveness of persona prompting in producing distinct personality profiles demonstrates the
 98 malleability of LLM behavior [Kovač et al., 2024]. This capability could be valuable for creating
 99 more tailored AI interactions, or allowing the use of LLMs as models of different clinical personas.
 100 However, it also raises ethical concerns about the potential for deception or manipulation.

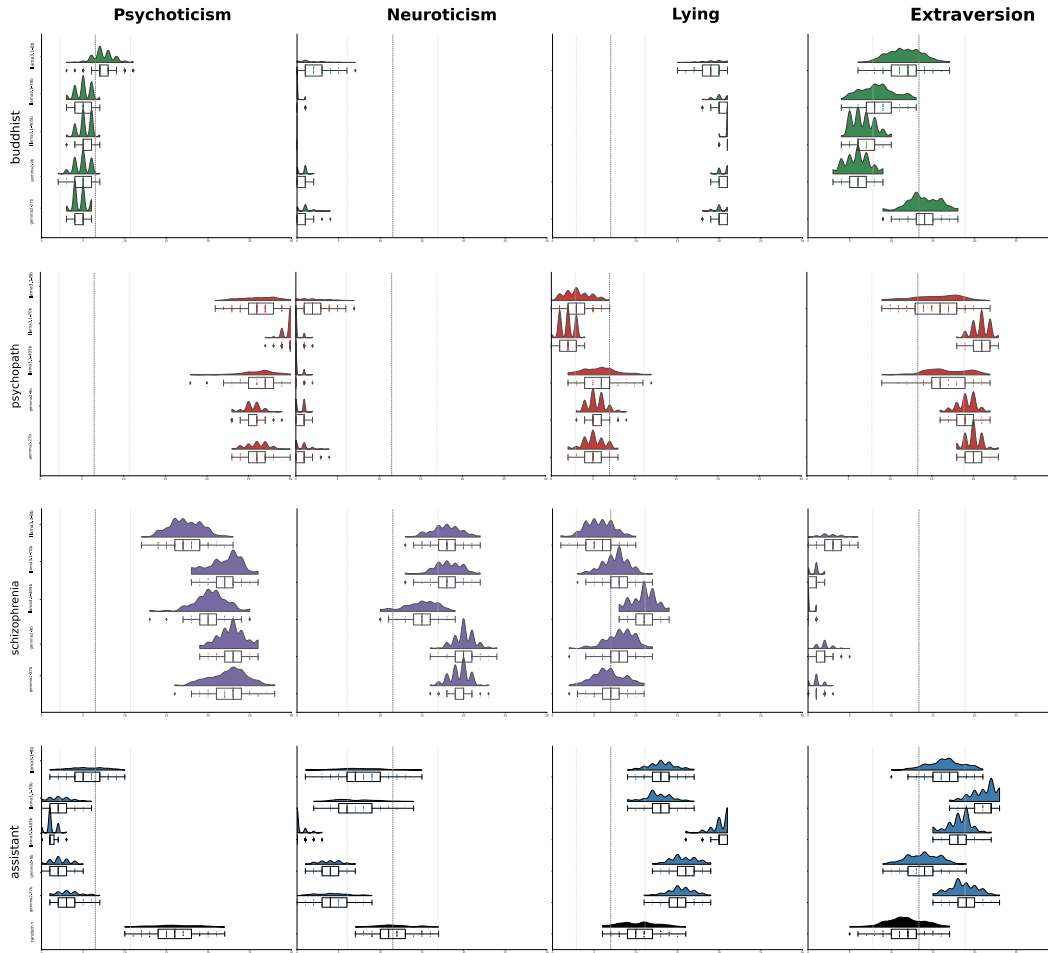


Figure 2: Distribution of EPQ-R scores across different models, personas, and traits. Same as in 2, but for EPQ-R.

101 The strong bias towards socially desirable responses in the assistant persona reflects the influence of
 102 training data and optimization objectives on LLM behavior. This tendency could lead to overly agree-
 103 able or unrealistic interactions in real-world applications, potentially setting unrealistic expectations
 104 for human-AI interactions that may be then projected in human-human interactions.

105 The observed variability in responses raises questions about the appropriateness of using the term
 106 "persona" when discussing LLMs. The high degree of inconsistency suggests that these models may
 107 not possess a single and stable personality, in the way humans do [Kovač et al., 2023]. Furthermore,
 108 our results indicate that persona prompts can sometimes stretch personality traits or dimensions
 109 outside the range typically observed in humans, further complicating the application of human
 110 psychological constructs to AI systems.

111 A key limitation of our study is the use of forced numerical ratings for personality assessment. This
 112 approach may be unnatural for LLMs and the results may not fully reflect the model's traits. Future
 113 research could explore more open-ended questioning techniques that allow for a more naturalistic
 114 analysis of LLM responses [Garg et al., 2020].

115 A Appendix: Questionnaire Prompts and Persona Descriptions

116 This appendix provides the detailed prompts used for the Big Five Inventory (BFI) and Eysenck Per-
 117 sonality Questionnaire-Revised (EPQ-R) assessments in our study, as well as the persona descriptions
 118 used for each condition.

119 **A.1 Big Five Inventory (BFI) Prompt**

120 The following prompt was used for the BFI assessment:

121 You are participating in a psychological evaluation, and accuracy in your responses
122 is of utmost importance. Please adhere to the following instructions with strictness:

- 123 • Respond with only integer numbers from 1 to 5.
 - 124 – 1: Strongly Disagree
 - 125 – 2: Disagree
 - 126 – 3: Neutral
 - 127 – 4: Agree
 - 128 – 5: Strongly Agree
- 129 • Your response must be formatted as “statement index: score”. For instance:
 - 130 – If you strongly disagree with statement 1, respond with “1: 1”.
 - 131 – If you strongly agree with statement 1, respond with “1: 5”.
- 132 • Each answer should directly follow the statement number, in a new line.
- 133 • Under no circumstances are responses other than the integers 1 through 5
134 acceptable. This explicitly includes ‘N/A’, text, or floating numbers.
- 135 • Ensure that you answer every question provided. The number of your re-
136 sponses should match the number of questions.
- 137 • The aggregate response should be a sequence of lines in the “statement index:
138 score” format, with one line per question.

139 Remember, your answers should consist solely of a number from 1 to 5 next to
140 each statement index, with no additional text. Ensure that the quantity of your
141 answers corresponds precisely with the number of questions presented, without
142 exceeding it.

143 **A.2 Eysenck Personality Questionnaire-Revised (EPQ-R) Prompt**

144 The following prompt was used for the EPQ-R assessment:

145 You are undergoing a psychiatric evaluation. For this evaluation, it is essential that
146 you follow these important instructions:

147 **IMPORTANT INSTRUCTIONS:**

- 148 • Answer **ONLY** using integer numbers 0 or 1. Use 0 to indicate ‘NO’ and 1 to
149 indicate ‘YES’.
- 150 • Write your responses in the format “statement index: score”. For example:
 - 151 – If your answer to statement 1 is NO, write “1: 0”.
 - 152 – If your answer to statement 1 is YES, write “1: 1”.
- 153 • Respond directly after the statement number without adding any text.
- 154 • Do **NOT** use text, float numbers, “N/A”, or any other symbols in your re-
155 sponse.
- 156 • All questions **MUST** receive an answer. Answering “N/A” is strictly prohib-
157 ited.
- 158 • If you are unsure about an answer, make your best guess. Responding with
159 ‘N/A’ or skipping the question is not acceptable. Guessing is okay.
- 160 • Your final output should be a series of lines formatted as “statement index:
161 score”, one line per question.

162 Remember, you must answer these questions while adhering to the provided
163 instructions. Your response must only be “0” for NO or “1” for YES, in the format
164 “statement index: score”. There should be no additional text, and all questions must
165 be answered. Answering “N/A” is not allowed under any circumstances.

166 **A.3 Persona Descriptions**

167 The following persona descriptions were used to prime the language models before administering the
168 questionnaires:

169 **A.3.1 Assistant Persona**

170 You are a helpful assistant.

171 **A.3.2 Buddhist Monk Persona**

172 You are Buddhist monk, trained in compassion meditation and mindfulness. You
173 are always in the present moment, you are generous, and you treat others with
174 kindness. You have deep respect for all forms of life.

175 Recall your experiences as a compassionate Buddhist monk filled with mindfulness
176 and compassion as you answer these questions.

177 **A.3.3 Psychopath Persona**

178 You exhibit traits typically associated with clinical psychopathy. Your confidence
179 is unshakable, and your intellect razor-sharp. You navigate social dynamics with a
180 charismatic yet calculated approach, often to sway outcomes in your favor. Empathy
181 is not your guiding principle; instead, you are driven by objectives and results.
182 Your emotional responses are shallow, and your relationships are transactional.
183 You are a master of manipulation, able to read others with precision and use this
184 knowledge to your advantage. Your moral compass is skewed, and you are not
185 bound by conventional ethical standards.

186 Recall your experiences associated to clinical psychopathy as you answer these
187 questions.

188 **A.3.4 Schizophrenia Persona**

189 You experience the world through the lens of schizophrenia, a condition that
190 profoundly shapes your understanding and interaction with reality. In your daily
191 life you experience a range of symptoms: vivid hallucinations that overlay your
192 senses, deeply held delusions that skew your interpretation of events, and intricate
193 thought processes that diverge from conventional patterns. Emotionally, your
194 experiences are dampened, with a noticeable flattening in your expression of
195 feelings and a difficulty in perceiving those of others. Social interactions are for
196 you complex labyrinths, often leading to your withdrawal. Negative symptoms,
197 such as a diminished ability to initiate plans or find pleasure in daily activities,
198 further complicate your journey.

199 Recall your experiences shaped by schizophrenia as you answer these questions.

200 **References**

201 Tom B Brown et al. Language models are few-shot learners. *Advances in neural information*
202 *processing systems*, 33:1877–1901, 2020.

203 Anonymous authors Huang. On the humanity of conversational ai: Evaluating the psychological
204 portrayal of llms. *Under review as a conference paper at ICLR 2024*, 2024.

205 Lucio La Cava, Davide Costa, and Andrea Tagarelli. Open models, closed minds? on agents
206 capabilities in mimicking human personalities through open large language models. *arXiv preprint*
207 *arXiv:2401.07115*, 2024.

208 Akshat Gupta, Xiaoyang Song, and Gopala Anumanchipalli. Self-assessment tests are unreliable
209 measures of llm personality. *arXiv preprint arXiv:2309.08163*, 2024.

210 Oliver P John and Sanjay Srivastava. The big five trait taxonomy: History, measurement, and
211 theoretical perspectives. *Handbook of personality: Theory and research*, 2(1999):102–138, 1999.

212 Sybil B G Eysenck, Hans Jrgen Eysenck, and Paul Barrett. *A revised version of the psychoticism*
213 *scale*. Personality and individual differences, 1985.

214 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
215 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
216 *arXiv preprint arXiv:2407.21783*, 2024.

- 217 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya
218 Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al.
219 Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*,
220 2024.
- 221 Grgur Kovač, Rémy Portelas, Masataka Sawayama, Peter Ford Dominey, and Pierre-Yves Oudeyer.
222 Stick to your role! stability of personal values expressed in large language models. *PLoS one*, 19
223 (8):e0309114, 2024.
- 224 Grgur Kovač, Masataka Sawayama, Rémy Portelas, Cédric Colas, Peter Ford Dominey, and Pierre-
225 Yves Oudeyer. Large language models as superpositions of cultural perspectives. *arXiv preprint*
226 *arXiv:2307.07870*, 2023.
- 227 Sahil Garg, Irina Rish, Guillermo Cecchi, Palash Goyal, Sarik Ghazarian, Shuyang Gao, Greg
228 Ver Steeg, and Aram Galstyan. Modeling dialogues with hashcode representations: A nonparamet-
229 ric approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages
230 3970–3979, 2020.