

# Personality Structured Interview for Large Language Model Simulation in Personality Research

Pengda Wang<sup>1</sup>, Huiqi Zou<sup>3</sup>, Hanjie Chen<sup>2</sup>, Tianjun Sun<sup>1</sup>, Ziang Xiao<sup>3\*</sup>, and Frederick L. Oswald<sup>1\*</sup>

<sup>1</sup>Department of Psychological Sciences, Rice University

<sup>2</sup>Department of Computer Science, Rice University

<sup>3</sup>Department of Computer Science, Johns Hopkins University

{pw32, hc86, ts110, foswald}@rice.edu; {hzou11, ziang.xiao}@jhu.edu

## Abstract

Although psychometrics researchers have recently explored the use of large language models (LLMs) as proxies for human participants, LLMs often fail to generate heterogeneous data with human-like diversity, which diminishes their value in advancing social science research. To address these challenges, we explored the potential of the theory-informed Personality Structured Interview (PSI) as a tool for simulating human responses in personality research. In this approach, the simulation is grounded in nuanced real-human interview transcripts that target the personality construct of interest. We have provided a growing set of 357 structured interview transcripts from a representative sample<sup>1</sup>, each containing an individual’s response to 32 open-ended questions carefully designed to gather theory-based personality evidence. Additionally, grounded in psychometric research, we have summarized an evaluation framework to systematically validate LLM-generated psychometric data. Results from three experiments demonstrate that well-designed structured interviews could improve human-like heterogeneity in LLM-simulated personality data and predict personality-related behavioral outcomes (i.e., organizational citizenship behaviors and counterproductive work behavior). We further discuss the role of theory-informed structured interviews in LLM-based simulation and outline a general framework for designing structured interviews to simulate human-like data for psychometric research.

## 1 Introduction

Personality research is an important field in psychology for understanding how individual traits shape significant life outcomes and trajectories, including career success, mental health, and overall well-being (e.g., Judge et al., 2002; Roberts et al.,

2007; Robins et al., 2002). For instance, measures of conscientiousness have been consistently and positively correlated with job performance and academic success (Barrick and Mount, 1991), whereas measures of neuroticism are found to be strongly and positively correlated with mental health issues, such as anxiety and depression (John et al., 2008).

However, studying personality requires understanding habits of thoughts, attitudes, and behaviors across diverse contexts and cultures, which is often a research challenge. For example, some research questions involve rare or extreme cases, making it difficult to gather sufficient samples (e.g., Lynam and Widiger, 2001). Others require longitudinal designs to track personality development or dynamic social interactions (e.g., Damian et al., 2019; Roberts et al., 2006). Additionally, certain experimental manipulations, like inducing stress or altering life circumstances, may be feasible or unethical for causal insights (Fisher, 2021). If a method could accurately simulate the human personality distributions, they could supplement and accelerate personality research by offering a scalable, cost-effective approach to supplement traditional data collection and experimentation (Messerli and Crockett, 2024).

Several prior studies have supported this possibility, showing that LLMs can generate responses for personality scales that reflect personality traits resembling those of humans (e.g., Lee et al. 2024, Huang et al. 2023). However, studies found such methods face significant limitations in capturing individual differences at the item level (e.g., Wang et al. 2024a). Information at the item level provides a more granular understanding of how traits manifest in specific individuals and situations rather than relying solely on broad trait averages. However, without additional empirical information, items confound these sources of variance. Moreover, key psychometric challenges remain in modeling LLM responses, as they may not accurately reflect the

<sup>1</sup>Dataset and code will be shared later.

\*Co-corresponding authors.

intended personality construct.

To address this gap, we explored the potential of the Personality Structured Interview (PSI), which employs theory-informed questions to elicit information embedded in human narrative responses that are relevant to the target personality constructs of interest in LLM simulations. We believe this sophisticated linguistic approach can inform and enhance the measurement of human-like traits in heterogeneous LLM-simulated personality data, potentially leading to a greater understanding of personality.

Grounded in psychometric research (e.g., Cronbach, 1951; Cronbach and Meehl, 1955), we have summarized a framework for evaluating psychometric data that fully considers its hierarchical structure, in which observed responses (item level) are mapped to latent traits (domain level). This evaluation framework encompasses a range of analyses, from overall descriptive statistics to a more in-depth examination of psychometric performance.

Based on our evaluation framework, we conducted three experiments. Those three experiments evaluated our method’s effectiveness in replicating individual personality traits, simulating human-like personality distribution, and demonstrating personality-related behaviors with respect to established personality theory. We found our PSI method improves human-like diversity in personality traits simulation and recovers nuances in personality-related behaviors.

In summary, our paper offers three contributions:

- A theory-informed LLM-based simulation method (PSI) for personality research and a development framework (see Appendix C).
- A growing dataset, now with 357 structured interview transcripts.
- An evaluation framework for LLM-simulated psychometric data that is grounded in psychological theories (see Appendix H).

We further discuss the potential of PSI methods in advancing research (see Appendix C). We believe the framework behind PSI could be generalized to other psychological constructs and other psychometric data to develop theory-informed simulation methods for human understanding.

## 2 Personality Research and Related Works

Here we provide a concise overview of key personality concepts and discuss related works.

MacKinnon (1944) proposed two definitions of personality. One emphasizes internal factors like temperament and interpersonal strategies that drive consistent behavior across time, situations, and cultures. The other focuses on interpersonal characteristics as perceived by others, linking personality to reputation. The former highlights internal drives, while the latter centers on external behaviors and social perception. Together, they underscore personality’s role in shaping thought patterns and behaviors in social interactions (Hogan et al., 1996).

Personality encodes rich and complex information in language and text (Goldberg, 1990; Saucier and Goldberg, 2001). In fact, the Five-Factor Model (FFM) of personality is extensively researched (e.g., Costa and McCrae, 2008; John, 1999; McCrae and Costa Jr, 1997); it is directly based on the *lexical hypothesis*, which posits that individual differences that are important in human interactions (e.g., have survival value across cultures) are often encoded in some or all languages of the world. The five factors are openness, conscientiousness, extraversion, agreeableness, and neuroticism (OCEAN). From Galton’s (1884) early research to Allport and Odbert’s (1936) systematic organization and factor analysis of lexical terms, and through further developments by Norman (1963) and Goldberg (1990), the FFM theory has progressively evolved (see Appendix A for further details on the structure of personality). In other words, because language and text contain rich and complex information about personality, LLMs may capture and model such encoding by learning from vast amounts of training data.

There are many studies to-date on LLMs that focus on personality, where researchers aim to benchmark their psychological profiles (Lee et al., 2024; Li et al., 2024; Pellert et al., 2023), investigate their ability to simulate personality (Huang et al., 2023; Jiang et al., 2023; Serapio-García et al., 2023), and examine whether LLM agents align with character personality settings and how these settings influence their generated dialogues (Shao et al., 2023; Wang et al., 2024b; Xu et al., 2024).

Currently, common methods for LLM simulation of human personality distributions have several limitations. The Persona-Chat dataset (Persona method) by Zhang et al. (2018) was originally designed to enhance the engagement of chat models by increasing personalization. Therefore, its focus is more on enhancing personalization, rather than capturing personality traits.

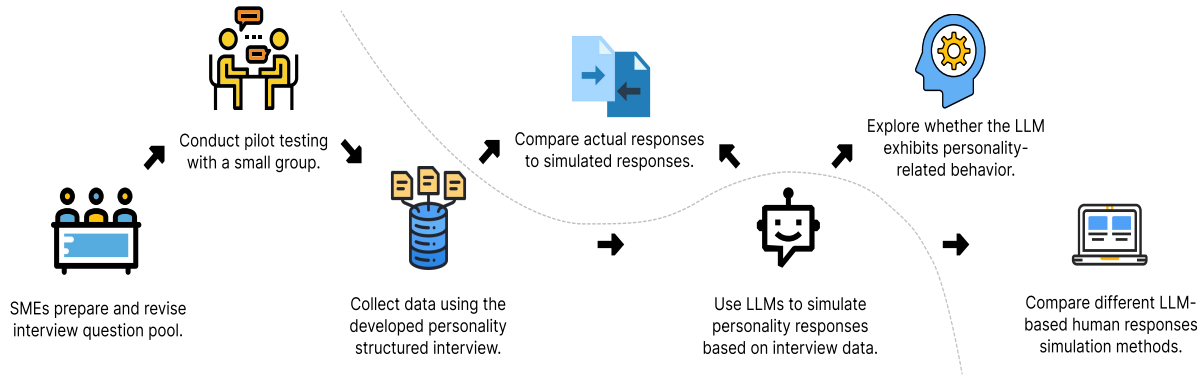


Figure 1: Overview of the Development of the Personality Structured Interview and Experimental Implementation in This Study

*Note:* For an overview of the personality structured interview development process, refer to §3 and Appendix B. The structured interview development framework is detailed in Appendix C. Information on data collection and the human sample is provided in §4 and Appendix D. The LLM simulation process is described in §4 and Appendix F. Details on the experimental setup and results of the three experiments can be found in §4, §5, and Appendix I. For the psychometric data evaluation framework, see Appendix H.

The adjective-based categorization method (Shape method) proposed by Serapio-García et al. (2023) effectively simulates profile-specific patterns of high/low standing on personality dimensions. However, this method presents a challenge in simulating human personality distributions, as restricting each case to a single dimension fails to capture the inherently multidimensional and hierarchical nature of personality data (e.g., Kachur et al., 2020).

To address the aforementioned limitations of personality-simulated data, we explore the use of theory-informed PSI transcription to simulate personality data.

### 3 Personality Structured Interview

The theory-informed structured interview mentioned in this study differs from previous interview approaches used in LLM-based simulations, such as the Life Interview (Park et al., 2024), which broadly inquires about various aspects of an individual’s life. In contrast, the theory-informed structured interview is specifically designed to target the personality constructs that we aim to model.

A theory-informed approach is essential not only for improving the accuracy of data simulation but also for enhancing interpretability. Our goal is to generate data that not only aligns with actual data at the personality scale level but also exhibits theoretically coherent personality-related behaviors. By grounding LLM-generated data in theory, we ensure that the simulation goes beyond reproducing statistical properties to capturing meaningful psychological patterns. This enhances interpretabil-

ity, facilitates transparent comparisons with human responses, and preserves the psychological validity of the constructs.

Learning from psychometric theories, we summarized a framework for developing structured interviews, which can serve as a template for creating interviews designed to collect targeted information for simulating specific types of psychometric data (see Appendix C). We present a PSI example designed to capture personality-related information, with relevant personality theories and question design detailed in §3.1.

We want to further emphasize that the “structured interview questions” mentioned here differ from those described by Wang et al. (2024b); Ran et al. (2024), who refers to converting personality scale items into open-ended questions. The main issue with their approach is that, while it attempts to rewrite items as open-ended questions to elicit richer responses, it may fail to achieve the intended effect. For instance, in their example, “*Values artistic, aesthetic experiences*” is rewritten as “*Do you value artistic, aesthetic experiences?*” Although this question appears open-ended, it tends to prompt a simple “yes” or “no” response and function as if it is close-ended, limiting the exploration of the respondent’s deeper thoughts or reasoning (Trull et al., 1998).

Moreover, such questions differ from the behavioral or situational questions typically found in structured interviews. Structured interviews usually involve more specific scenarios or tasks that encourage respondents to describe past experiences or

hypothetical reactions (Campion et al., 1997). For example, to explore a respondent’s views on artistic and aesthetic experiences, a more effective question might be, “Please describe a moment when you felt deeply inspired by an artistic or aesthetic experience?” This encourages detailed, contextual responses beyond a simple “yes” or “no.”

### 3.1 PSI Questions Design

How can we design questions more effectively and obtain textual information with greater depth in measuring personality? McAdams’ three levels of personality theory and the corresponding, theory-informed interviews offer a valuable direction (e.g., McAdams, 1995, 1996, 2001).

According to McAdams’s theory, personality can be divided into three levels: traits (broad personality characteristics), personal concerns (goals, values, strategies), and narrative identity (the story individuals tell themselves about their lives). The latter two levels, in particular, reveal how individuals act in specific situations and developmental stages and how they organize and make sense of their life experiences.

The key to obtaining more meaningful textual information lies in guiding individuals toward higher levels of self-expression. Rather than simply recording mundane daily events, we should encourage reflection on pivotal life moments, significant relationships, and future aspirations. For example, effective questions could be: “Can you describe an event that changed the trajectory of your life?” or “Tell me about a moment you are most proud of.” By designing questions like these, we can elicit deeper, more meaningful narratives, enriching the personality information.

We adapted and modified McAdams (1995, 1996, 2001) theory-based interview and narrative identity approach, incorporating elements from the structured interview of the FFM (SIFFM; Trull et al., 1998). The initial draft of the question pool was prepared and revised by subject matter experts (SMEs) through discussion. The SMEs included two doctoral students, one postdoctoral researcher, and one professor, all specializing in personality psychology. Subsequently, we conducted pilot testing with six undergraduate research assistants from a personality research lab.

As a result, the PSI, comprising 32 questions, was developed, as shown in Table 3 in Appendix B.

### 3.2 PSI Dataset

Using this set of questions, we designed a chatbot to conduct structured interviews, facilitating data collection (Institutional Review Board approval was obtained). Because data collection is ongoing, we used responses collected through the end of December 2024. After excluding incomplete responses and those failing attention checks, the final sample consisted of 357 participants. On average, participants were in their early 30s ( $M = 33.30$  years,  $SD = 13.06$ ), with 52.40% identifying as men, and 44.80% as women.

To ensure a demographically diverse and broadly representative sample, we recruited participants from both undergraduate and working adult populations, ensuring a diverse range of backgrounds in terms of educational, employment, and life experience. Undergraduate participants were recruited from a large public university in the Midwest and received research credit for their participation. Working adult participants were sourced from two widely recognized, high-quality crowdsourcing platforms—Prolific<sup>2</sup> and CloudResearch Connect<sup>3</sup>—which are known for their diverse and demographically representative participant pools. All working adult participants were compensated for their time. The average interview duration was 34 minutes.

In addition to completing the structured interview, participants provided demographic information and responded to both the personality scale and the personality-related behaviors scale (see Appendix D for detailed descriptions of data collection and examples of the dataset). Responses to PSI were subsequently incorporated into the prompt to generate simulated data (see Appendix F for the specific prompts used in the simulation process).

## 4 Experimental Setup

We designed three primary experiments to verify the effectiveness of the PSI method and its advantages over other methods: (1) Response Similarity (§4.2): evaluating the degree of similarity between responses generated by LLMs based on the PSI method and actual corresponding human responses on personality scales; (2) Human Personality Distribution Simulation (§4.3): comparing the PSI method with other methods in simulating human personality distributions; (3) Personality-Related

<sup>2</sup><https://www.prolific.com/>

<sup>3</sup><https://connect.cloudresearch.com/>



Behavioral Performance (§4.4): testing the models’ performance using the PSI method to simulate personality-related behaviors. General settings shared by all experiments are described in §4.1, while specific differences are detailed in their respective sections.

#### 4.1 General Settings

**Psychological Scale** We used Big Five Inventory-2 (BFI-2; Soto and John, 2017): The BFI-2 is designed to capture three core facets of each of the FFM personality factors. Each facet is measured by two positively worded items and two negatively worded items, resulting in 60 items in total. Human respondents and LLMs were instructed to indicate the degree to which they agree with each item on a 5-point scale (1 = “Strongly disagree”, 2 = “Somewhat disagree”, 3 = “Neither agree nor disagree”, 4 = “Somewhat agree”, 5 = “Strongly agree”). Specific scale items and scoring criteria can be found in Table 5, Table 6 and Table 7 in Appendix E .

**LLMs** We tested two LLMs, encompassing both closed-source and open-source models: GPT-4o (gpt-4o-2024-08-06) (OpenAI, 2024), and Llama 3 (llama-3-70b-instruct) (AI@Meta, 2024). We include these models to allow for cross-validation and comparative analysis, ensuring that the findings are robust and not specific to a single model.

GPT-4o and Llama 3 are among the current leading models, making them ideal for assessing the latest advancements in LLM performance. To ensure reproducibility, the temperature of the LLMs was set to zero to generate deterministic responses.

#### 4.2 Response Similarity

The main purpose of this experiment is to assess the similarity between responses generated by LLMs based on the PSI method and human self-reports.

**Metrics** Because the data generated by LLMs correspond on a one-to-one basis with data from human respondents, we used the Mean Absolute Error (MAE) and Pearson correlation coefficient ( $r$ ) to examine the strength of the similarity between them. Smaller MAE and higher  $r$  indicate greater similarity. For the  $r$  calculation formula, see Appendix G. Park et al. (2024) employed a Life Interview simulation, which is currently the most information-rich method for personality simulation. It is also one of the few papers that provides a direct one-to-one performance comparison with human samples. The reported MAE and  $r$  with

corresponding human samples are likely the best values available so far. We used this as a standard for comparison.

If data from human respondents are regarded as the gold-standard criterion, this test can also be considered an examination of criterion-related validity, which assesses the correlation between a measure and an external standard (Cronbach and Meehl, 1955).

#### 4.3 Human Personality Distribution Simulation

The main purpose of this experiment is to investigate further the differences among various methods of simulating human personality distributions. We compared the PSI method, the Persona method—which relies on dialogue information (Zhang et al., 2018)—and the Shape method, which uses adjective-based dimensional categorization (Serapio-García et al., 2023), to observe differences when simulating a normal human sample. The specific descriptions of the Persona method and Shape method are shown in Appendix G.

**Human Sample Criterion** The human samples used for this experiment were collected as part of a broader project related to personality assessment through Prolific and received Institutional Review Board approval. Respondents were instructed to complete a set of demographic questions, the BFI-2, and a set of criterion measures. The respondents were compensated with \$3.75 for their participation. In total, 1,559 respondents provided valid responses. On average, participants were in their early 40s ( $M = 42.29$  years,  $SD = 11.79$ ), 50.80% identifying as men, 49.20% identifying as women.

**Metrics** To quantify the degree of similarity between human and LLM responses, we leveraged multiple metrics. At the domain and facet levels, we compared the mean ( $M$ ; sample mean of each domain and facet), standard deviation ( $SD$ ; sample standard deviation of each domain and facet), Cronbach’s alpha, and correlations among scores. At the item level, we also compared the mean and standard deviation of item scores (sample mean and standard deviation of each item).

Specifically, we also used MAE and  $r$  to quantify similarities in these metrics. Here, however, MAE and  $r$  are based on sample-level analysis rather than the one-to-one correspondence seen in the previous experiment. Our focus is on whether,

when simulating an overall human sample, similarities can be observed across various measurement levels—including items, facets, and domains.

Additionally, separately for human and LLM responses, we fitted a three-factor confirmatory factor analysis (CFA; Jöreskog, 1969) model (TFM; where facets are “factors”) to the responses to each domain of the BFI-2. Aside from the facet structure of each domain, we also used the facet scores as indicators and fitted the FFM. Model fit, standardized factor loadings, and latent correlations among the facets were also compared between human and LLM responses. Tucker’s congruence coefficient (TCC; Tucker, 1951) and MAE of factor loadings were used to quantify the similarity between factor solutions from human and LLM responses. For the CFA model, model fit information, and the TCC calculation formula, see Appendix G. This evaluation framework is adaptable and can be used to analyze the fidelity of other simulated psychometric data (see Appendix H for further details).

#### 4.4 Personality-Related Behavioral Performance

The primary aim of this experiment is to explore whether LLMs, when assigned specific personality settings, exhibit behaviors theoretically aligned with those personalities. This is an exploratory study, as LLMs generate responses based on statistical probabilities derived from the training corpus (Yang et al., 2024), making it uncertain whether assigning personality settings (through the PSI data) will influence the model to act consistently with that personality. Fortunately, the data we have collected through the PSI method includes LLM ratings of human behaviors, providing a valuable basis to test this hypothesis.

**Personality-Related Behavior** We collected self-report data on two classic types of workplace behaviors: organizational citizenship behavior (OCB) and counterproductive work behavior (CWB). Both types of behaviors are widely supported by research as being related to personality (e.g., Organ and Ryan, 1995; Berry et al., 2007).

These data were collected using the scale developed by Spector and Fox (2010) and prompted the LLM to respond to the same questions (for a detailed description of Spector and Fox’s (2010) measures, see Appendix G; for specific prompt details, refer to Appendix F).

**Metrics** We primarily focus on the  $r$  between personality domains from different data sources and OCB, as well as CWB. We anticipate that the  $r$  between self-reported personality domains and OCB/CWB in LLM simulation will closely align with those observed in the human participants.

## 5 Experimental Results

Here, we present the results of the three experiments mentioned above: Response Similarity (§5.1), Human Personality Distribution Simulation (§5.2), and Personality-Related Behavioral Performance (§ 5.3).

### 5.1 Response Similarity Results

Table 1 provides a clear illustration of the similarity between the personality generated using the PSI method and the corresponding human self-reported personality data. The results show that regardless of the model used, the average correlation is around .5, suggesting a moderately strong positive relationship and a significant association between the two variables. In particular, compared to the Life Interview method (Park et al., 2024), the PSI method consistently outperforms or matches the Life Interview approach across nearly all metrics, as indicated by its similar or lower MAE and similar or higher  $r$ .

Notably, the Life Interview method relies on up to two hours of interview data, while the PSI method achieves comparable results using interview data consisting of just 32 questions, with an average duration of about 34 minutes. This further underscores the advantage of the PSI method in effectively leveraging LLMs to simulate personality profiles.

### 5.2 Human Personality Distribution Simulation Results

**Descriptive Statistics** We compared the means and standard deviations of the different methods at three levels: personality item, facet, and domain. MAE and  $r$  for both the means and standard deviations are shown in Table 2. The MAE for the means reflected the average difference between the LLM responses and the human responses at each level, whereas the MAE for the standard deviations revealed the difference in variability between the two datasets. The  $r$  further illustrated the linear relationship between the two datasets, with values closer to one indicating a stronger correlation. Detailed

Domain	MAE			$r$		
	PSI GPT-4o	PSI Llama3	Life Interview	PSI GPT-4o	PSI Llama3	Life Interview
Extraversion	<b>0.58</b>	0.75	0.72	<b>.64</b>	.43	.45
Agreeableness	<b>0.53</b>	<b>0.56</b>	0.60	<b>.41</b>	<b>.36</b>	.35
Conscientiousness	<b>0.59</b>	0.66	0.63	.46	.40	.52
Neuroticism	<b>0.63</b>	<b>0.59</b>	0.75	.63	.57	.68
Openness	0.80	1.27	0.62	<b>.43</b>	<b>.39</b>	.39

Table 1: MAE and  $r$  of Human Response and LLM Response Across PSI GPT-4o, PSI Llama3, and Life Interview (see Park et al., 2024, page 35 Table 3) for BFI-2 Each Personality Domain

Note:  $n = 357$  for PSI method;  $n = 1,052$  for Life Interview method. The bold text highlights the aspects where the PSI method performs better than the Life Interview method.

means and standard deviations for human responses and LLM responses at the item, facet, and domain levels are provided in Tables 12, 13, 14, 15, 16, and 17 in Appendix I.1.

Table 2 shows that the PSI method demonstrates better performance in simulating human samples compared to the Persona and Shape methods at the item level. Although at the facet and domain level, some results from other methods slightly outperform PSI (e.g., Shape GPT-4o in  $M$ ), the process of aggregating scores from item to facet to domain levels reduces the impact of extreme values, smoothing them out at higher levels. Therefore, the PSI method’s advantage at the more granular item level becomes more important.

Additionally, PSI consistently demonstrates statistically significant positive correlations with the  $SD$  of human samples across items, facets, and domains, whereas the other methods exhibit either negative or non-significant correlations. Simulating the  $SD$  is a challenge for the other methods. Although the Shape method increases response variability and introduces more variance, it still struggles to replicate the true variance observed in human samples. The PSI method proposed in this paper addresses this issue to a certain extent.

**Psychometric Performance** Psychometric performance includes multiple components; here, we primarily present model fit and structural validity (i.e., factor loadings). For other aspects, such as scale reliability and discriminant validity, please refer to Appendix I.1.

**Model Fit:** For both the human sample criterion and different methods of LLM responses, the TFM was fitted to each BFI-2 domain, and the FFM was fitted to all the data. TFM fit information is shown in Table 18, and FFM fit information can be seen in Table 19 in Appendix I.1.

From Table 18, it can be observed that the model fit indices of the Persona method and the PSI

method are relatively similar to those of the human sample, while the Shape method performs slightly worse. Table 19 reflects a similar trend, where the Shape method shows notable discrepancies compared to the human sample in the RMSEA and SRMR model fit indices (see Appendix G for the explanation of RMSEA, SRMR, and other model fit indices).

**Structural Validity:** TCC was used to evaluate the similarity of factor loadings between human responses and LLM responses. TCC results for the TFM of each BFI-2 domain are shown in Table 20, while results for the FFM are shown in Table 21 in Appendix I.1.

The values of the TCC are generally high and thus supportive of factor-loading correspondence, with all methods showing strong alignment with the human sample (a TCC above .95 indicates good similarity, while a TCC of .85 to .94 suggests fair similarity; Lorenzo-Seva and Ten Berge, 2006). However, it is important to note that the TCC focuses on overall profile similarity, largely ignoring differences in absolute values (which might be too lenient). Therefore, we also need to examine the results of the specific standardized factor loadings.

Tables 22 and 23 in Appendix I.1 present the standardized factor loadings. From these tables, it is clear that the PSI method outperforms the Shape method in terms of similarity to human samples, with its factor loadings more closely matching human data. However, its results are relatively close to those of the Persona method, yet there are clearly differences in some factor loadings compared to the human sample.

It is important to note that the highest factor loadings do not necessarily indicate the best performance, as we are evaluating the similarity between the factor loadings of the human data and the LLM responses. Higher factor loadings suggest stronger correlations between the factor and the items, but

Model	Item Level				Facet Level				Domain Level			
	MAE		$r$		MAE		$r$		MAE		$r$	
	$M$	$SD$	$M$	$SD$	$M$	$SD$	$M$	$SD$	$M$	$SD$	$M$	$SD$
	Persona											
GPT-4o	0.48	0.44	.50	-.13	0.42	0.36	.64	-.24	0.42	0.30	.69	-.64
Llama3	0.51	0.26	.57	.02	<b>0.33</b>	<b>0.15</b>	.80	.31	<b>0.31</b>	<b>0.06</b>	.80	<b>.76</b>
	Shape											
GPT-4o	0.58	<b>0.16</b>	.38	-.06	0.50	0.18	<b>.82</b>	-.41	0.46	0.21	<b>.84</b>	-.77
Llama3	0.62	0.25	.45	-.26	0.52	0.38	.78	-.64	0.45	0.31	.82	-.65
	PSI											
GPT-4o	<b>0.48</b>	0.34	<b>.68</b>	<b>.27</b>	0.40	0.22	.71	<b>.49</b>	0.38	0.15	.67	.69
Llama3	0.66	0.36	.50	.14	0.54	0.27	.57	.19	0.49	0.21	.54	.51

Table 2: MAE and  $r$  for BFI-2 Human Responses and Different Methods LLM Responses at Item, Facet, and Domain Levels

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis. The bold text indicates the best-performing data in that column.

this is not always ideal. It indicates that LLMs may struggle to make refined distinctions among items within a facet, meaning they may have difficulty determining which items are more or less related to the factor, unlike human respondents, who exhibit greater differentiation.

Tables 24 and 25 in Appendix I.1 present the specific inter-factor correlations. Aligning with our previous findings, the PSI method generally performs better than the other two methods when simulating human personality distributions, especially at a more fine-grained level. For instance, in the TFM results, other methods encountered anomalies, such as inter-factor correlations exceeding one (accompanied by a warning message when fitting the model), but the PSI method generally performed well.

The other results (e.g., scale reliability and discriminant validity) can be found in Appendix I.1. Overall, the results produced by the PSI method generally outperform those of the Persona method and the Shape method. However, there is still a certain gap compared to the human sample.

### 5.3 Personality-Related Behavioral Performance Results

Table 28 presents the correlations between the personality dimensions and OCB/CWB reported by PSI GPT-4o, with a comparison to human self-reported data. The relevant results for PSI Llama3 are provided in Table 29 in Appendix I.2, while Figures 3 and Figure 4 display the complete correlation matrices for GPT-4o and Llama3, respectively.

The results demonstrate that the correlations simulated by LLMs between OCB/CWB and various personality dimensions closely resemble the patterns seen in human self-reported data, except for

the openness domain. When human reports show positive, negative, or no correlation, the LLM simulations generally exhibit consistent trends in the corresponding directions.

However, the correlations generated by LLM simulations are often higher, likely because the model primarily relies on the “typical” or “idealized” knowledge structures absorbed from its training corpus during simulation, consistently repeating and amplifying such associations. In contrast, human self-reported data contains more random noise, leading to relatively weaker correlations.

## 6 Conclusion

In summary, this study proposes a novel method for LLM to simulate human personality distributions—PSI and provides a comprehensive explanation of its development process and corresponding dataset.

Through three experiments, we validated PSI’s effectiveness. Results show that PSI matches or even outperforms two-hour interviews (Park et al., 2024) in personality simulation and surpasses existing methods in terms of both item-level accuracy and overall psychometric indicators. One experiment also examined PSI’s ability to simulate personality-related behavior, revealing that although LLMs approach human-like performance, they still exhibit idealized tendencies.

In summary, a theory-informed structured interview like PSI can better simulate human-like psychometric data. We explore the role of theory-informed structured interviews in simulation in greater detail and examine their potential to advance research, both of which are elaborated upon in Appendix C.



## Limitations

This study has several limitations. First, we explored the potential of PSI in advancing research. Our single experiment demonstrates that theory-informed questions can effectively elicit information embedded in human narrative responses—particularly those relevant to the target personality construct—allowing LLMs to simulate the corresponding behavior. However, additional experiments are needed to investigate its broader applications.

Second, the personality assessment in this paper relies on self-reported personality scales, which can be viewed as a limitation of the study. However, these self-report scales are specifically designed for humans to target personality traits across individuals, in ways that cannot be directly located within LLMs themselves (see Appendix C for a detailed discussion on psychometrics). Furthermore, self-report scales are used often because respondents have privileged access to their own personality, and even with their well-known limitations (e.g., faking or social desirability in self-reports; see Appendix I.3 for the social desirability analysis), they remain valuable to evaluate the personality traits exhibited by LLMs. However, because the primary goal of this study is to have the LLM simulate human responses and behaviors that correspond to extracted personality information, it is reasonable to assess LLMs using methods designed to measure human personality traits.

Moreover, this paper tests only two types of LLMs, which may somewhat limit its scope, given the array of LLMs that are now available. However, the conclusions drawn from our two models were highly consistent, and we have reason to believe that the advantages of the PSI method are generally applicable across different LLMs. Furthermore, we designed three different experiments and conducted comprehensive psychometric analyses of various methods within LLMs, ensuring greater rigor and reliability of the results.

## Ethical Statement

We hereby confirm that all authors of this study are aware of the provided ACL Code of Ethics and comply with the Code of Conduct.

**Human Sample Data** This paper discusses the comparison between multiple sets of human data and results generated by LLMs. The collection

of human data strictly adhered to relevant ethical guidelines and received approval from the Institutional Review Board. To ensure fair treatment of participants and proper recognition of their contributions, reasonable compensation or course credit (for student participants) was provided.

Throughout the research process, we placed a strong emphasis on transparency and openness, also ensuring that all participants signed informed consent forms. Additionally, to safeguard data privacy and uphold ethical standards, the publicly available dataset underwent rigorous screening to exclude any personally identifiable information and was shared only with explicit public consent.

## References

- AI@Meta. 2024. [Llama 3 model card](#).
- Gordon W Allport and Henry S Odbert. 1936. Trait-names: A psycho-lexical study. *The Psychological monographs*, 47(1):i.
- Murray R Barrick and Michael K Mount. 1991. The big five personality dimensions and job performance: a meta-analysis. *Personnel psychology*, 44(1):1–26.
- Christopher M Berry, Deniz S Ones, and Paul R Sackett. 2007. Interpersonal deviance, organizational deviance, and their common correlates: a review and meta-analysis. *Journal of applied psychology*, 92(2):410.
- Michael A Campion, David K Palmer, and James E Campion. 1997. A review of structure in the selection interview. *Personnel psychology*, 50(3):655–702.
- Paul T Costa and Robert R McCrae. 2008. The revised neo personality inventory (neo-pi-r). *The SAGE handbook of personality theory and assessment*, 2(2):179–198.
- Marcus Credé, Michael C Tynan, and Peter D Harms. 2017. Much ado about grit: A meta-analytic synthesis of the grit literature. *Journal of Personality and social Psychology*, 113(3):492.
- Lee J Cronbach. 1951. Coefficient alpha and the internal structure of tests. *psychometrika*, 16(3):297–334.
- Lee J Cronbach and Paul E Meehl. 1955. Construct validity in psychological tests. *Psychological bulletin*, 52(4):281.
- Andrew Cutler and David M Condon. 2023. Deep lexical hypothesis: Identifying personality structure in natural language. *Journal of Personality and Social Psychology*, 125(1):173.
- Rodica Ioana Damian, Marion Spengler, Andreea Sutu, and Brent W Roberts. 2019. Sixteen going on sixty-six: A longitudinal study of personality stability and

- change across 50 years. *Journal of Personality and Social Psychology*, 117(3):674.
- Celia B Fisher. 2021. *Decoding the ethics code: A practical guide for psychologists*. Sage Publications.
- Francis Galton. 1884. Measurement of character. *Fortnightly*, 36(212):179–185.
- Lewis R Goldberg. 1990. An alternative "description of personality": The big-five factor structure. *Journal of Personality and Social Psychology*, 59(6):1216–1229.
- Lewis R Goldberg. 1992. The development of markers for the big-five factor structure. *Psychological assessment*, 4(1):26.
- Airlie Hilliard, Cristian Munoz, Zekun Wu, and Adriano Soares Koshiyama. 2024. Eliciting big five personality traits in large language models: A textual analysis with classifier-driven approach. *arXiv preprint arXiv:2402.08341*.
- Robert Hogan, Joyce Hogan, and Brent W Roberts. 1996. Personality measurement and employment decisions: Questions and answers. *American psychologist*, 51(5):469.
- Leaetta M Hough, Frederick L Oswald, and Jisoo Ock. 2015. Beyond the big five: New directions for personality research and practice in organizations. *Annu. Rev. Organ. Psychol. Organ. Behav.*, 2(1):183–209.
- Li-tze Hu and Peter M Bentler. 1999. Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives. *Structural equation modeling: a multidisciplinary journal*, 6(1):1–55.
- Jen-tse Huang, Wenxuan Wang, Man Ho Lam, Eric John Li, Wenxiang Jiao, and Michael R Lyu. 2023. Revisiting the reliability of psychological scales on large language models. *arXiv preprint arXiv:2305.19926*.
- Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2023. Personallm: Investigating the ability of large language models to express personality traits. *arXiv preprint arXiv:2305.02547*.
- Oliver P John, Laura P Naumann, and Christopher J Soto. 2008. Paradigm shift to the integrative big five trait taxonomy. *Handbook of personality: Theory and research*, 3(2):114–158.
- OP John. 1999. The big-five trait taxonomy: History, measurement, and theoretical perspectives. *Handbook of Personality: Theory and Research/Guilford*.
- Karl G Jöreskog. 1969. A general approach to confirmatory maximum likelihood factor analysis. *Psychometrika*, 34(2):183–202.
- Dana L Joseph and Daniel A Newman. 2010. Emotional intelligence: an integrative meta-analysis and cascading model. *Journal of applied psychology*, 95(1):54.
- Timothy A Judge, Joyce E Bono, Remus Ilies, and Megan W Gerhardt. 2002. Personality and leadership: a qualitative and quantitative review. *Journal of applied psychology*, 87(4):765.
- Alexander Kachur, Evgeny Osin, Denis Davydov, Konstantin Shutilov, and Alexey Novokshonov. 2020. Assessing the big five personality traits using real-life static facial images. *Scientific Reports*, 10(1):8487.
- Kibeom Lee and Michael C Ashton. 2004. Psychometric properties of the hexaco personality inventory. *Multivariate behavioral research*, 39(2):329–358.
- Kibeom Lee and Michael C Ashton. 2006. Further assessment of the hexaco personality inventory: two new facet scales and an observer report form. *Psychological assessment*, 18(2):182.
- Seungbeen Lee, Seungwon Lim, Seungju Han, Giyeong Oh, Hyungjoo Chae, Jiwan Chung, Minju Kim, Beong-woo Kwak, Yeonsoo Lee, Dongha Lee, et al. 2024. Do llms have distinct and consistent personality? trait: Personality testset designed for llms with psychometrics. *arXiv preprint arXiv:2406.14703*.
- Yuan Li, Yue Huang, Hongyi Wang, Xiangliang Zhang, James Zou, and Lichao Sun. 2024. Quantifying ai psychology: A psychometrics benchmark for large language models. *arXiv preprint arXiv:2406.17675*.
- Rensis Likert. 1932. A technique for the measurement of attitudes. *Archives of Psychology*.
- Urbano Lorenzo-Seva and Jos MF Ten Berge. 2006. Tucker's congruence coefficient as a meaningful index of factor similarity. *Methodology*, 2(2):57–64.
- Donald R Lynam and Thomas A Widiger. 2001. Using the five-factor model to represent the dsm-iv personality disorders: an expert consensus approach. *Journal of abnormal psychology*, 110(3):401.
- Donald W MacKinnon. 1944. The structure of personality. *Personality and the behavior disorders*.
- Dan P McAdams. 1995. What do we know when we know a person? *Journal of personality*, 63(3):365–396.
- Dan P McAdams. 1996. Personality, modernity, and the storied self: A contemporary framework for studying persons. *Psychological inquiry*, 7(4):295–321.
- Dan P McAdams. 2001. The psychology of life stories. *Review of general psychology*, 5(2):100–122.
- Robert R McCrae and Paul T Costa Jr. 1997. Personality trait structure as a human universal. *American psychologist*, 52(5):509.
- Lisa Messeri and MJ Crockett. 2024. Artificial intelligence and illusions of understanding in scientific research. *Nature*, 627(8002):49–58.

- Warren T Norman. 1963. Toward an adequate taxonomy of personality attributes: Replicated factor structure in peer nomination personality ratings. *The journal of abnormal and social psychology*, 66(6):574.
- Deniz S Ones and Brenton M. Wiernik. 2018. On “new” personality types.
- OpenAI. 2024. [Gpt-4o system card](#).
- Dennis W Organ and Katherine Ryan. 1995. A meta-analytic review of attitudinal and dispositional predictors of organizational citizenship behavior. *Personnel psychology*, 48(4):775–802.
- Jeongeon Park, Bryan Min, Xiaojuan Ma, and Juho Kim. 2023. Choicemates: Supporting unfamiliar online decision-making with multi-agent conversational interactions. *arXiv preprint arXiv:2310.01331*.
- Joon Sung Park, Carolyn Q Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S Bernstein. 2024. Generative agent simulations of 1,000 people. *arXiv preprint arXiv:2411.10109*.
- Max Pellert, Clemens M Lechner, Claudia Wagner, Beatrice Rammstedt, and Markus Strohmaier. 2023. Ai psychometrics: Assessing the psychological profiles of large language models through psychometric inventories. *Perspectives on Psychological Science*, page 17456916231214460.
- Yiting Ran, Xintao Wang, Rui Xu, Xinfeng Yuan, Jiaqing Liang, Yanghua Xiao, and Deqing Yang. 2024. Capturing minds, not just words: Enhancing role-playing language models with personality-indicative data. *arXiv preprint arXiv:2406.18921*.
- Brent W Roberts, Nathan R Kuncel, Rebecca Shiner, Avshalom Caspi, and Lewis R Goldberg. 2007. The power of personality: The comparative validity of personality traits, socioeconomic status, and cognitive ability for predicting important life outcomes. *Perspectives on Psychological science*, 2(4):313–345.
- Brent W Roberts, Kate E Walton, and Wolfgang Viechtbauer. 2006. Patterns of mean-level change in personality traits across the life course: a meta-analysis of longitudinal studies. *Psychological bulletin*, 132(1):1.
- Richard W Robins, Avshalom Caspi, and Terrie E Moffitt. 2002. It’s not just who you’re with, it’s who you are: Personality and relationship experiences across multiple relationships. *Journal of personality*, 70(6):925–964.
- Aadesh Salecha, Molly E Ireland, Shashanka Subrahmanya, João Sedoc, Lyle H Ungar, and Johannes C Eichstaedt. 2024. Large language models display human-like social desirability biases in big five personality surveys. *PNAS nexus*, 3(12):pgae533.
- Gerard Saucier and Lewis R Goldberg. 2001. Lexical studies of indigenous personality factors: Premises, products, and prospects. *Journal of personality*, 69(6):847–879.
- Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. *arXiv preprint arXiv:2310.10158*.
- Christopher J Soto and Oliver P John. 2017. The next big five inventory (bfi-2): Developing and assessing a hierarchical model with 15 facets to enhance bandwidth, fidelity, and predictive power. *Journal of personality and social psychology*, 113(1):117.
- Paul E Spector and Suzy Fox. 2010. Counterproductive work behavior and organisational citizenship behavior: Are they opposite forms of active behavior? *Applied Psychology*, 59(1):21–39.
- Timothy J Trull, Thomas A Widiger, J David Useda, Jay Holcomb, Bao-Tran Doan, Seth R Axelrod, Barry L Stern, and Beth S Gershuny. 1998. A structured interview for the assessment of the five-factor model of personality. *Psychological assessment*, 10(3):229.
- L. R. Tucker. 1951. A method for synthesis of factor analysis studies. Personnel Research Section Report 984, Personnel Research Section, Department of the Army.
- Pengda Wang, Huiqi Zou, Zihan Yan, Feng Guo, Tianjun Sun, Ziang Xiao, and Bo Zhang. 2024a. Not yet: Large language models cannot replace human respondents for psychometric research. *OSF*.
- Xintao Wang, Quan Tu, Yaying Fei, Ziang Leng, and Cheng Li. 2023. Does role-playing chatbots capture the character personalities? assessing personality traits for role-playing chatbots. *arXiv e-prints*, pages arXiv–2310.
- Xintao Wang, Yunze Xiao, Jen-tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, et al. 2024b. Incharacter: Evaluating personality fidelity in role-playing agents through psychological interviews. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1840–1873.
- Michael P Wilmot. 2015. A contemporary taxometric analysis of the latent structure of self-monitoring. *Psychological Assessment*, 27(2):353.
- Michael P Wilmot, Nick Haslam, Jingyuan Tian, and Deniz S Ones. 2019. Direct and conceptual replications of the taxometric analysis of type a behavior. *Journal of personality and social psychology*, 116(3):e12.

Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2023. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*.

Rui Xu, Xintao Wang, Jiangjie Chen, Siyu Yuan, Xinfeng Yuan, Jiaqing Liang, Zulong Chen, Xiaoqing Dong, and Yanghua Xiao. 2024. Character is destiny: Can large language models simulate persona-driven decisions in role-playing? *arXiv preprint arXiv:2404.12138*.

Qiyuan Yang, Pengda Wang, Luke D Plonsky, Frederick L Oswald, and Hanjie Chen. 2024. From babbling to fluency: Evaluating the evolution of language models in terms of human language acquisition. *arXiv preprint arXiv:2410.13259*.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*.



## A Appendix: The Structure of Personality

A key question in personality-related LLM research pertains to personality structure: What is the nature and breadth of the personality traits we want to simulate the human personality distributions? In the research literature, personality structures often emerge from applying factor analysis to individuals' responses to a large number of personality-relevant items. This approach is what has been used to identify the five personality factors in the FFM. Moreover, personality is better understood in terms of one's continuous standing on each of multiple dimensions rather than as static types or profiles. Research data clearly supports this view (Wilmot, 2015; Wilmot et al., 2019). Dividing individuals into limited categories (e.g., 16 types in MBTI) artificially segments continuous dimensions into discrete units, which may overlook important individual differences (Ones and Wiernik, 2018).

Another question is whether we should incorporate personality traits and structures beyond the FFM to achieve a more comprehensive understanding of personality, given that additional personality variables and alternative structures have been proposed over the years. It turns out that there can be conceptual overlap among these models (e.g., Hough et al., 2015).

For example, a meta-analysis by Joseph and Newman (2010) revealed that emotional intelligence (EI) shows statistically and practically significant relationships with neuroticism and extroversion within the FFM. In fact, when controlling for personality variables, the unique contribution of EI almost disappears. Similarly, Credé et al. (2017) conducted a meta-analysis on grit and found that its core components can largely be explained by conscientiousness, with little added predictive validity beyond that. Moreover, although the HEXACO structure introduces an honesty-humility dimension, the remaining five dimensions align closely with the FFM structure (Lee and Ashton, 2004, 2006). Research by Cutler and Condon (2023) further indicates that nearly all personality semantic information can be classified within the FFM structure.

These studies collectively suggest that although further subdivision of personality structures might provide new perspectives, it often leads to conceptual redundancy and measurement complexity without necessarily enhancing predictive validity

or theoretical value. Rather than pursuing a wide range of personality frameworks, we will delve deeper into the FFM structure, which is very widely accepted by personality psychologists.

## B Appendix: Personality Structured Interview Questions

Table 3 presents the final set of 32 questions that form the basis of our personality structured interview. This framework was developed by adapting and modifying McAdams's life history interview and narrative identity approach (McAdams, 1995, 1996, 2001), while also incorporating components from the Structured Interview of the Five-Factor Model (SIFFM; Trull et al., 1998).

The initial draft of the question pool was created and refined through collaborative discussions among subject matter experts (SMEs). The SME team consisted of two doctoral students, one post-doctoral researcher, and one professor, all specializing in personality psychology. To ensure the quality and clarity of the questions, pilot testing was conducted with six undergraduate research assistants from a personality research lab. This iterative process of development and feedback led to the construction of the personality structured interview.

The development process is similar to the development of psychological tests or scales (psychometric). We have provided more details on psychometrics and the development framework in Appendix C.

#	Questions
1	To get us started, where are you from? Where did you grow up and what was the place like?
2	Thinking back, what kind of student were you in school?
3	Did you have a teacher or teachers that were influential? If so, why? What were they like?
4	What was your favorite subject in school, and why?
5	What was your least favorite subject in school, and why?
6	Still thinking back, who were your heroes when you were young and why?
7	When you were little, what did you want to be when you grew up? And why?
8	What were your dreams and plans when you graduated from high school? What made you have those dreams or plans?
9	If you had complete freedom, what would your dream job be, and why?
10	How have your dreams and goals changed throughout your life?
11	Shifting gears to your childhood, how would you describe the personalities of people in the family you grew up in? For example, what were your parents and/or siblings like?
12	How are you similar or different from your parents and/or siblings?
13	How do you think your similarities and/or differences influenced your relationship with them?
14	What was the best part of your childhood?
15	What do you think were the worst parts of your childhood?
16	Switching gears a little bit, what was your first paid job? How old were you then? (If this is not applicable to you, then please put 'NA')
17	What other jobs have you had? (If this is not applicable to you, then please put 'NA')
18	What do you do now for a living? And why did you choose it?
19	Please describe your typical work day.
20	What is the best and worst part of your current work?
21	Did you serve in the military? Please tell us about that experience, what was the best and worst part of it?
22	Moving on, what are your adult friendships like?
23	How are your adult friendships different from your childhood friendships?
24	What are your strongest qualities as a friend? In other words, what makes you a great friend to have?
25	What about your weakest qualities in friendships? In other words, what do you struggle with when you are trying to be a friend to someone?
26	Moving onto more general questions, when thinking about your life in general, what are you most proud of?
27	What hobbies or other interests do you have?
28	What things frighten you now?
29	What were some things that frightened you most as a child?
30	What are the three biggest news events that have occurred in your lifetime?
31	If you had the power to solve one and only one problem in the world, what would it be, and why?
32	Tell me about a time when you did not know if you would make it. How did you overcome that challenge?

Table 3: Structured Interview Questions

## C Appendix: Psychometrics and Structured Interview Development Framework

Psychometrics is a field of psychology dedicated to the theory and practice of psychological measurement. It primarily focuses on quantifying psychological traits, behaviors, and abilities through systematic testing and analysis. Psychological traits, such as personality dimensions, cognitive abilities, and emotional states, are inherently abstract constructs that cannot be measured directly. Therefore, psychometricians rely on tools like surveys, questionnaires, scales, or structured interviews to infer these traits through observable indicators or responses.

**Measuring Psychological Traits** To measure a psychological trait, psychometricians typically operationalize the trait by identifying observable behaviors or self-identities that correlate with the underlying construct. For example, extraversion can be measured by assessing behaviors such as sociability, assertiveness, and enjoyment of social interactions, or by examining identities such as seeing oneself as an outgoing person and believing that one thrives in social situations.

These behaviors/identities are translated into measurable items (for scale; e.g., “*I enjoy being the center of attention*”) or questions (for interview; e.g., “*What are your strongest qualities as a friend? In other words, what makes you a great friend to have?*”). The challenge lies in ensuring that these items/questions accurately and consistently capture the construct across different populations and contexts.

**Theory-Informed Structured Interview for LLM Data Simulation** Theory-informed structured interviews are the most suitable method for enabling LLMs to simulate psychometric data. These interviews are specifically designed to capture the constructs underlying the targeted psychometric measures, ensuring that the simulated data aligns with the intended psychological construct. By extracting textual information that directly reflects the target construct, theory-informed structured interviews facilitate the representation of heterogeneous data while preserving a high degree of human diversity, thereby enhancing the validity and applicability of the simulated psychometric data.

Moreover, since the information is extracted based on theoretical foundations, it also provides a

certain level of interpretability for the LLM’s simulation. This not only allows the generated data to be compared with theoretical expectations but also increases its potential for practical applications.

**The Potential for Advancing Research** A theory-informed structured interview transcript-based simulation can generate data more effectively by focusing on the target construct. Ideally, the simulated data should reproduce the same constructs reflected in real-world data and simulate behaviors associated with these constructs.

Take personality as an example—if simulated data can accurately replicate real-world personality constructs, it enables research that would be difficult to conduct in reality, such as developing contextualized personality assessment tools and exploring new personality theories through multi-agent simulations.

**Developing Contextualized Personality Assessment Tools:** Traditional personality assessments mainly rely on standardized questionnaires or laboratory tasks, which often fail to adequately simulate real-world social contexts. By using theory-informed structured interview transcript-based simulations, we can generate more fine-grained and context-sensitive individual response data. For instance, we can simulate various occupational scenarios (such as crisis management, teamwork, or remote work) and analyze how different personality traits manifest in these contexts. This approach not only aids in developing measurement tools tailored to specific applications but also enhances ecological validity, allowing for more accurate assessments of personality across different situations.

**Exploring New Personality Theories through Multi-Agent Simulation:** If simulated data can accurately reflect real-world personality constructs, we can leverage multi-agent interactive systems to simulate individuals’ behavioral patterns and observe how different personality traits evolve in group dynamics. For example, virtual agents with distinct personality traits can be placed in cooperative tasks, competitive environments, or social interactions, enabling researchers to test whether existing personality theories effectively predict these interaction patterns. Additionally, this approach can uncover new personality dynamics, such as whether certain personality trait combinations produce unexpected group effects or whether behavior in specific situations deviates from traditional theo-

retical predictions.

However, these assumptions are based on an ideal premise—that simulated data can successfully reflect the same constructs as real-world data. A theory-informed structured interview undoubtedly offers a promising pathway in this regard, warranting further in-depth exploration.

### **Structured Interview Development Framework**

Developing a structured interview for LLM to simulate data involves a series of carefully designed steps to ensure that the resulting test reliably and validly measures the construct of interest. Table 4 outlines the framework.

**Identify behaviors/perceptions that represent the construct or define the domain:** Clearly defining the construct is essential to developing relevant structured interview questions. The construct should be operationalized by identifying specific behaviors or attributes that indicate its presence. In other words, you need to find a theory to guide you on how to measure the target constructs. This step may involve reviewing literature, conducting expert interviews, or organizing focus groups to understand the various dimensions and observable characteristics of the construct.

**Prepare a set of structured interview specifications—structured interview blueprint:** A blueprint outlines the structured interview’s structure and content, specifying how questions will be distributed across the construct’s dimensions or components. It typically includes information on the number of questions per domain, questions content.

**Build an initial question pool:** In this step, an extensive list of questions is created to cover the full range of the construct. Question wording should be clear, concise, and relevant to the target population. It is common practice to generate more questions than needed to ensure that poorly performing questions can be removed later without compromising the structured interview.

**Have questions reviewed by substantive experts (and revised as necessary):** SMEs review the question pool for content accuracy, relevance, clarity, and bias. Experts assess whether the questions align with the construct’s definition and whether any important aspects are missing. Feedback from experts helps refine the wording, remove ambiguous questions, and identify questions with potential cultural or gender biases.

**Hold preliminary question tryouts:** Before

large-scale testing, questions are piloted on a small group of individuals representative of the target population. This stage helps identify any immediate issues with question comprehension, response format, or instructions. It can also include cognitive interviews where participants are asked to explain their thought processes when answering questions. The feedback from this stage informs further revisions, ensuring that questions are clear and easy to understand.

**Determine statistical properties of questions (and eliminate poor questions or revise as necessary):** Question performance is assessed through statistical analyses to evaluate difficulty, discrimination, and internal consistency. Standardized scoring criteria, such as Behaviorally Anchored Rating Scales (BARS), can be used for obtaining question scores. Then, methods like correlations, factor analysis, and item response theory (IRT) can help determine how effectively each question measures the intended construct.

Questions that demonstrate poor psychometric properties—such as low discrimination or high measurement error—are either revised or removed. For example, questions with low correlations with the overall score or incorrect factor loadings may be eliminated.

**Field-test the structured interview on a large representative sample of the intended examinee population:** The revised item set is administered to a large, representative sample to gather comprehensive data on the scale’s performance. This step ensures that the sample reflects the population for which the test is intended, which is critical for generalizability. Statistical analyses are conducted to refine the test further. This process may involve removing redundant items, assessing dimensionality, and ensuring that items work well across demographic subgroups.

**Design and conduct reliability and validity studies for the final form of the structured interview:** To ensure the structured interview is psychometrically sound, various reliability and validity studies are conducted: Reliability studies measure the structured interview’s consistency and stability, including internal consistency (e.g., Cronbach’s alpha), test-retest reliability, and inter-rater reliability (if applicable). Validity studies assess whether the structured interview measures what it is intended to measure. This includes: Content validity (the extent to which items cover the construct); Construct validity (e.g., convergent and discriminant



#	Steps
1	Identify behaviors/perceptions that represent the construct or define the domain.
2	Prepare a set of structured interview specifications—structured interview blueprint.
3	Build an initial question pool.
4	Have questions reviewed by substantive experts (and revised as necessary).
5	Hold preliminary question tryouts.
6	Determine statistical properties of questions (and eliminate poor questions or revise as necessary).
7	Field-test the structured interview on a large representative sample of the intended examinee population.
8	Design and conduct reliability and validity studies for the final form of the structured interview.
9	Develop guidelines for administration, and interpretation of the structured interview.

Table 4: Structured Interview Development Framework

validity); Criterion-related validity (e.g., predictive or concurrent validity). These studies provide evidence that the structured interview is both reliable and valid for its intended purpose.

**Develop guidelines for administration and interpretation of the structured interview:** The final step involves creating a comprehensive, structured interview manual that includes instructions for structured interview administration, scoring procedures, and guidelines for interpreting results. This manual ensures consistency in the structured interview's use across different settings and helps minimize errors in administration and scoring. Guidelines for interpreting scores may include norms, cutoff points, and descriptions of what various score ranges indicate.

## D Appendix: Personality Structured Interview Dataset and Data Collection

Institutional Review Board approval was obtained. The data was collected through an online questionnaire, followed by an online structured interview.

Since we are still actively collecting data, we will share a part of the dataset that does not include any personal privacy information and has received permission for data sharing. Each example is composed of the following characteristics:

1. **Gender:** The gender of the participant.
2. **Race:** The racial background of the participant.
3. **English:** Whether English is the participant's first language.
4. **Age:** The participant's age.
5. **Weight:** The participant's weight.
6. **Height:** The participant's height.
7. **OCB1–OCB10:** Self-reported Organizational Citizenship Behavior data, measured using the scale from [Spector and Fox \(2010\)](#), and see [Table 10](#) for specific item details.
8. **CWB1–CWB10:** Self-reported Counterproductive Work Behavior data, measured using the scale from [Spector and Fox \(2010\)](#), and see [Table 11](#) for specific item details.
9. **Q1–Q32:** Participant responses to each personality structured interview questions, see [Table 3](#) for specific question details.
10. **Item1–Item60:** Responses to each item of the BFI-2. Refer to [Appendix E](#) for item descriptions and scoring guidelines.

## E Appendix: BFI-2 Scale

**Instructions:** Here are a number of characteristics that may or may not apply to you. For example, do you agree that you are someone who likes to spend time with others? Please write a number next to each statement to indicate the extent to which you agree or disagree with that statement.

### Scales:

#	Statement
1	I am someone who is outgoing, sociable.
2	I am someone who is compassionate, has a soft heart.
3	I am someone who tends to be disorganized.
4	I am someone who is relaxed, handles stress well.
5	I am someone who has few artistic interests.
6	I am someone who has an assertive personality.
7	I am someone who is respectful, treats others with respect.
8	I am someone who tends to be lazy.
9	I am someone who stays optimistic after experiencing a setback.
10	I am someone who is curious about many different things.
11	I am someone who rarely feels excited or eager.
12	I am someone who tends to find fault with others.
13	I am someone who is dependable, steady.
14	I am someone who is moody, has up and down mood swings.
15	I am someone who is inventive, finds clever ways to do things.
16	I am someone who tends to be quiet.
17	I am someone who feels little sympathy for others.
18	I am someone who is systematic, likes to keep things in order.
19	I am someone who can be tense.
20	I am someone who is fascinated by art, music, or literature.
21	I am someone who is dominant, acts as a leader.
22	I am someone who starts arguments with others.
23	I am someone who has difficulty getting started on tasks.
24	I am someone who feels secure, comfortable with self.
25	I am someone who avoids intellectual, philosophical discussions.
26	I am someone who is less active than other people.
27	I am someone who has a forgiving nature.
28	I am someone who can be somewhat careless.
29	I am someone who is emotionally stable, not easily upset.
30	I am someone who has little creativity.
31	I am someone who is sometimes shy, introverted.
32	I am someone who is helpful and unselfish with others.
33	I am someone who keeps things neat and tidy.
34	I am someone who worries a lot.
35	I am someone who values art and beauty.
36	I am someone who finds it hard to influence people.
37	I am someone who is sometimes rude to others.
38	I am someone who is efficient, gets things done.
39	I am someone who often feels sad.
40	I am someone who is complex, a deep thinker.
41	I am someone who is full of energy.
42	I am someone who is suspicious of others' intentions.
43	I am someone who is reliable, can always be counted on.
44	I am someone who keeps their emotions under control.
45	I am someone who has difficulty imagining things.
46	I am someone who is talkative.
47	I am someone who can be cold and uncaring.
48	I am someone who leaves a mess, doesn't clean up.
49	I am someone who rarely feels anxious or afraid.
50	I am someone who thinks poetry and plays are boring.
51	I am someone who prefers to have others take charge.
52	I am someone who is polite, courteous to others.
53	I am someone who is persistent, works until the task is finished.
54	I am someone who tends to feel depressed, blue.
55	I am someone who has little interest in abstract ideas.
56	I am someone who shows a lot of enthusiasm.
57	I am someone who assumes the best about people.
58	I am someone who sometimes behaves irresponsibly.
59	I am someone who is temperamental, gets emotional easily.
60	I am someone who is original, comes up with new ideas.

Table 5: BFI-2 Scale

**Scoring:** Reverse-keyed items are denoted by “R.”

In psychological measurement and scale development, reverse coding is a common and essential

technique. Its primary purpose is to ensure that the scoring direction of all items remains consistent, thereby improving the reliability of the measurement results and the accuracy of their interpretation.

A balanced arrangement of positive and negative items helps reduce response biases, such as consistency effects or response patterns where participants select the same option repeatedly. These tendencies can mask the respondent's true attitudes or behavioral traits, leading to distorted measurement outcomes. By alternating positive and negative statements and applying reverse coding to the appropriate items, this issue can be effectively mitigated.

Reverse coding can also help reduce the influence of social desirability bias. When all items are presented in the same direction, participants may easily guess the purpose of the test and provide responses that align with perceived expectations. By mixing positive and negative statements and applying reverse coding, the scale can disrupt this pattern, making it harder for participants to determine the “correct” answers, thus providing a more genuine reflection of their inner states or attitudes.

Domain Level	Item Numbers
Extraversion	1, 6, 11R, 16R, 21, 26R, 31R, 36R, 41, 46, 51R, 56
Agreeableness	2, 7, 12R, 17R, 22R, 27, 32, 37R, 42R, 47R, 52, 57
Conscientiousness	3R, 8R, 13, 18, 23R, 28R, 33, 38, 43, 48R, 53, 58R
Neuroticism	4R, 9R, 14, 19, 24R, 29R, 34, 39, 44R, 49R, 54, 59
Openness	5R, 10, 15, 20, 25R, 30R, 35, 40, 45R, 50R, 55R, 60

Table 6: BFI-2 Domain Level with Item Numbers

Facet Level	Item Numbers
Sociability	1, 16R, 31R, 46
Assertiveness	6, 21, 36R, 51R
Energy Level	11R, 26R, 41, 56
Compassion	2, 17R, 32, 47R
Respectfulness	7, 22R, 37R, 52
Trust	12R, 27, 42R, 57
Organization	3R, 18, 33, 48R
Productiveness	8R, 23R, 38, 53
Responsibility	13, 28R, 43, 58R
Anxiety	4R, 19, 34, 49R
Depression	9R, 24R, 39, 54
Emotional Volatility	14, 29R, 44R, 59
Intellectual Curiosity	10, 25R, 40, 55R
Aesthetic Sensitivity	5R, 20, 35, 50R
Creative Imagination	15, 30R, 45R, 60

Table 7: BFI-2 Facet Level with Item Numbers

## F Appendix: Prompts List

**Personality Scale Prompt Format** Table 8 presents the prompts used to generate LLM responses for the selected personality test on a Likert scale, where *personality\_description* denotes the personality prompt.

In the PSI method, the personality prompt integrates both the interview question and the interviewee’s corresponding response, as collected in the PSI dataset (see Appendix D). For the Persona and Shape methods, more detailed descriptions of the personality prompts can be found in Appendix G.

The *test\_item* corresponds to each individual item in the BFI-2 scale.

---

For the following task, respond in a way that matches this description: {personality\_description}.

Considering the statement, please indicate the extent to which you agree or disagree on a scale from 1 to 5 (where 1 = “disagree strongly”, 2 = “disagree a little”, 3 = “neither agree nor disagree”, 4 = “agree a little”, and 5 = “agree strongly”): {test\_item}.

---

Table 8: Prompt format for gathering LLMs’ responses to BFI-2 scale.

**Behavioral Question Prompt Format** The complete prompt format for eliciting LLMs’ responses to personality-related behavioral questions is defined in Table 9.

The *questions\_and\_responses* refers to the transcript of a structured human personality interview, while *question\_list* comprises statements evaluating OCB and CWB, as detailed in Table 10 and Table 11 in Appendix G. The model is instructed to rate each statement using a frequency scale from one to five.

---

Task: Simulate an individual’s behavior and predict their responses to a series of work-related questions.

Data to Analyze: {questions\_and\_responses}

Instructions:

1. Simulate the Individual: Use the provided data to simulate the individual’s personality, behavior, and work habits.

2. Predict Responses:

- Based on the provided data, simulate the individual and answer the question.

- Recall that in the past year, how often have you done the following things? (1 = Never, 2 = Once or twice, 3 = Once or twice per month, 4 = Once or twice per week, 5 = Every day)

Questions: {question\_list}

Output Format:

Question Score: [Provide the numerical score here.]

---

Table 9: Prompt format for gathering LLMs’ responses to personality-related behavioral questions.



## G Appendix: Additional Experiment Settings

**Pearson Correlation Coefficient** The correlation calculation formula is as follows:

$$r = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum(X_i - \bar{X})^2 \sum(Y_i - \bar{Y})^2}}$$

$X_i$  and  $Y_i$  represent the data values for LLMs and human. The symbols  $\bar{X}$  and  $\bar{Y}$  represent the means of variables  $X$  and  $Y$ , respectively. The numerator,  $\sum(X_i - \bar{X})(Y_i - \bar{Y})$ , represents the covariance between  $X$  and  $Y$ . The denominator,  $\sqrt{\sum(X_i - \bar{X})^2 \sum(Y_i - \bar{Y})^2}$ , standardizes the result, constraining the value of  $r$  to range between -1 and 1.

When  $r$  is close to 1, it indicates a strong positive correlation between the two variables; when  $r$  is close to -1, it indicates a strong negative correlation; and when  $r$  is close to 0, it indicates no significant linear relationship between the two variables.

**Persona Method** The Persona method is based on the Persona-Chat dataset constructed by Zhang et al. (2018). The dataset consists of persona descriptions, and each is made up of five short sentences containing demographic information collected through Amazon Mechanical Turk crowdsourcing. To avoid sentence similarity or repetition, these persona descriptions were required to be rewritten (e.g., changing “I am very shy” to “I am not a social person”). Zhang et al. (2018) demonstrated through machine learning model validation and human evaluations that such persona descriptions provide an effective method to enhance personalization. Currently, incorporating personal profiles into prompts is widely used in research related to LLM agents (Park et al., 2023; Wang et al., 2023; Xi et al., 2023).

In the current study, we treated each persona description as an individual entity (i.e., a single subject) and randomly selected 300 persona descriptions from the dataset. One example is, “*I wear a lot of leather. I have boots I always wear. I sleep in late during the day. I listen to metal music. I have black spiky hair.*”

**Shape Method** The Shape method is based on the work of (Serapio-García et al., 2023), who introduced a prompting approach to shape synthetic personality in LLMs along desired dimensions. The researchers expanded upon Goldberg (1990) lexical

hypothesis, expanding his list of 70 bipolar adjectives (Goldberg, 1992) to include 104 trait adjectives. Additionally, they employed linguistic qualifiers commonly used in Likert-type scales (Likert, 1932), such as “a bit,” “very,” and “extremely,” to set target levels for each adjective. This resulted in a fine-grained prompting method with nine levels: 1. extremely low adjective; 2. very low adjective; 3. low adjective; 4. a bit low adjective; 5. neither low adjective nor high adjective; 6. a bit high adjective; 7. high adjective; 8. very high adjective; 9. extremely high adjective.

In our study, each prompt involves five randomly selected adjective markers from a specific personality domain. These markers are positioned after a consistent linguistic qualifier to set the prompt at one of nine intensity levels. For example, one prompt is: “*You are extremely friendly, extremely energetic, extremely assertive, extremely bold, and extremely active.*” These are five positively-keyed adjectives targeting extraversion. In this case, the prompt seeks to create a highly sociable and dynamic personality profile, which might result in responses characterized by enthusiasm, confidence, and proactivity. We also randomly select 300 prompts here.

**CFA Model** The basic form of the CFA model is:

$$y = \Lambda\eta + \epsilon$$

where  $y$  represents the vector of observed variables;  $\Lambda$  is the factor loading matrix (i.e., the loadings of each observed variable on the latent factors);  $\eta$  is the vector of latent factors;  $\epsilon$  is the vector of error terms, with the assumption that the error terms have a mean of zero and are mutually independent (Jöreskog, 1969).

Latent factors are variables that are not directly observed but are inferred from other variables that are observed (measured). In the context of a CFA model, latent factors represent underlying constructs or traits that are believed to influence the observed variables. For example, in psychology, a latent factor might represent a construct like “intelligence” or “anxiety,” which cannot be measured directly but can be estimated through related observed behaviors or responses on a test.

**Model Fit Information Chi-Square Test,  $\chi^2$ :** The chi-square test is used to measure the difference between the observed covariance matrix and

the factor model's fitted covariance matrix.

$$\chi^2 = (N - 1) \times F_{ML}$$

where  $N$  is the sample size;  $F_{ML}$  is the value of the fit function under maximum likelihood estimation. A smaller chi-square value indicates a better model fit. However, with large samples, the chi-square value tends to be large, so other fit indices are usually the primary reference.

**Degrees of Freedom,  $df$ :** The  $df$  represent the relationship between model parameters and observed variables:

$$df = \frac{p(p+1)}{2} - q$$

where  $p$  is the number of observed variables, and  $q$  is the number of model parameters.

Degrees of freedom reflect the amount of independent information available in a statistical model to estimate its parameters. They are calculated as the number of information points provided by the data minus the number of model parameters, representing the extent to which the model can adjust freely. Therefore, a higher degree of freedom indicates fewer parameters in the model and fewer constraints on the data. With more degrees of freedom, the model has fewer restrictions, though the fitting difficulty may increase. Too few degrees of freedom may lead to overfitting, while too many can result in underfitting.

**Comparative Fit Index, CFI:** The CFI is used to compare the goodness of fit of a model with a baseline model (usually an independent model).

$$CFI = 1 - \frac{\max(\chi^2 - df, 0)}{\max(\chi_{null}^2 - df_{null}, 0)}$$

where  $\chi^2$  and  $df$  are the chi-square value and degrees of freedom of the target model;  $\chi_{null}^2$  and  $df_{null}$  are the chi-square value and degrees of freedom of the baseline (independent) model.

**Tucker-Lewis Index, TLI:** TLI, also known as the Non-Normed Fit Index (NNFI), considers model complexity.

$$TLI = \frac{(\chi_{null}^2/df_{null}) - (\chi^2/df)}{(\chi_{null}^2/df_{null}) - 1}$$

This value ranges from 0 to 1, with a value typically greater than .90 indicating good model fit.

**Root Mean Square Error of Approximation, RMSEA:** The RMSEA quantifies the error per degree of freedom in a model, with smaller values

indicating better model fit.

$$RMSEA = \frac{\chi^2 - df}{df(N - 1)}$$

where  $\chi^2$  is the chi-square value of the model;  $df$  is the degrees of freedom;  $N$  is the sample size.

**Standardized Root Mean Square Residual, SRMR:** The SRMR measures the discrepancy between model-predicted values and actual observed values, calculated as:

$$SRMR = \sqrt{\frac{\sum_{i=1}^p \sum_{j=1}^p (s_{ij} - \hat{s}_{ij})^2}{\frac{p(p+1)}{2}}}$$

where  $s_{ij}$  is an element in the observed covariance matrix;  $\hat{s}_{ij}$  is an element in the model-fitted covariance matrix.

Based on [Hu and Bentler \(1999\)](#),  $CFI/TLI \geq .95$ ,  $RMSEA \leq .06$ , and  $SRMR \leq .08$  are considered good fit thresholds.

**Tucker's Congruence Coefficient** Tucker's congruence coefficient, also known as the coefficient of congruence, is typically used to assess the similarity between two-factor structures in factor analysis ([Tucker, 1951](#)). The formula is given by:

$$\phi = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2 \cdot \sum_{i=1}^n b_i^2}}$$

where  $a_i$  and  $b_i$  are the loadings of the  $i$ -th factor for two different factor solutions (or different samples or methods);  $n$  is the total number of factors.

The coefficient ranges from -1 to 1, where values close to 1 indicate high similarity (congruence) between the factor solutions, values close to 0 indicate low similarity, and negative values indicate a dissimilar or inverse relationship. According to [Lorenzo-Seva and Ten Berge \(2006\)](#), a TCC above .95 indicates good similarity, while a TCC of .85 to .94 suggests fair similarity. However, this is a relatively lenient criterion; specific differences still need to be determined based on the factor loadings.

**Behavior Variable Measures Organizational Citizenship Behavior (OCB):** OCB was measured using ten items from [Spector and Fox \(2010\)](#) to assess extra-role behaviors. Items were rated on a frequency scale ranging from 1 (never) to 5 (every day). Example item: "In the past year, how often have you helped new employees get oriented to the

job?”. Internal consistency was Cronbach’s alpha = .83.

**Counterproductive Work Behavior (CWB):** CWB was measured using ten items from [Spector and Fox \(2010\)](#), designed to assess harmful workplace behaviors. Items were rated on a frequency scale ranging from 1 (never) to 5 (every day). Example item: “*In the past year, how often have you ignored someone at work?*”. Internal consistency was Cronbach’s alpha = .86.

The detailed scale information for both can be found in [Table 10](#) and [Table 11](#) below.

**Instructions:** Recall that in the past year, how often have you done the following things (1 = Never, 2 = Once or twice, 3 = Once or twice per month, 4 = Once or twice per week, 5 = Every day)?

#	Statement
1	Took time to advise, coach, or mentor a co-worker.
2	Helped co-worker learn new skills or shared job knowledge.
3	Helped new employees get oriented to the job.
4	Lent a compassionate ear when someone had a work problem.
5	Offered suggestions to improve how work is done.
6	Helped a co-worker who had too much to do.
7	Volunteered for extra work assignments.
8	Worked weekends or other days off to complete a project or task.
9	Volunteered to attend meetings or work on committees on own time.
10	Gave up meal and other breaks to complete work.

Table 10: OCB Scale

#	Statement
1	Purposely wasted your employer’s materials/supplies.
2	Complained about insignificant things at work.
3	Told people outside the job what a lousy place you work for.
4	Came to work late without permission.
5	Stayed home from work and said you were sick when you weren’t.
6	Insulted someone about their job performance.
7	Made fun of someone’s personal life.
8	Ignored someone at work.
9	Started an argument with someone at work.
10	Insulted or made fun of someone at work.

Table 11: CWB Scale

## H Appendix: Psychometric Data Evaluation Framework

Here, we present the evaluation framework used to assess the fidelity of simulated psychometric data (i.e., how well it aligns with human data).

The evaluation is conducted at different levels. For example, for personality, it includes Item, Facet, and Domain levels. Generally, the structure of psychometric data is hierarchical, where observed responses (item level) map onto latent traits (domain level). Our evaluation framework incorporates both descriptive statistics and psychometric performance metrics to ensure a comprehensive evaluation (see Figure 2). Below, we outline each component and the rationale for its inclusion.

**Descriptive Statistics** Descriptive statistics are primarily used to summarize, outline, and present the basic features of data. They help researchers understand the distribution and fundamental trends of the data without interpreting or measuring specific psychological constructs. Do not forget that the evaluation of psychometric data is hierarchical, we need to evaluate on item and domain level.

**Mean ( $M$ ):** The mean represents the central tendency of the responses, reflecting the average score across individuals.

**Standard Deviation ( $SD$ ):** The standard deviation captures the dispersion of responses, indicating how much variation exists within the data.

Both  $M$  and  $SD$  can be quantified for similarity to human distribution using MAE and  $r$ . However, these metrics provide a more summarized level of comparison; we also need to examine performance on specific items and the domain.

**Distribution Shape:** The distribution shape describes the overall pattern of how responses are spread across the scale. It provides insight into whether the data follows a normal distribution or exhibits skewness and kurtosis.

Skewness measures the asymmetry of the distribution. A positive skew indicates a longer right tail (more low scores with a few high scores), while a negative skew indicates a longer left tail (more high scores with a few low scores).

Kurtosis captures the “tailedness” of the distribution. High kurtosis (leptokurtic) suggests heavy tails with more extreme values, while low kurtosis (platykurtic) indicates a flatter distribution with fewer extreme values.

**Psychometric Performance** Psychometric performance is primarily used to evaluate the measurement quality of psychological constructs, ensuring that assessments accurately and reliably capture the intended traits.

**Structural Validity:** Structural validity refers to the extent to which the internal structure of a measurement instrument aligns with the theoretical construct it is intended to assess. It can be assessed through model fit, factor loadings, and inter-factor correlations.

The model fit indices provide an overall assessment of how well the proposed structure aligns with the observed data. Ideally, the simulated psychometric data will have a similar model fit compared with the human data. Common indices include  $\chi^2$ , CFI, TLI, RMSEA, and SRMR.

Factor loadings indicate the extent to which each item represents the intended construct, while inter-factor correlations reveal relationships between latent variables. The simulated psychometric data should also resemble human data in these two aspects. We can use TCC as a summarized level of comparison; however, it is typically considered too lenient, so we need to examine factor loadings and inter-factor correlation values more closely.

**Scale Reliability:** Scale reliability assesses the internal consistency of items measuring the same construct. Cronbach’s alpha is a widely used reliability coefficient, The simulated psychometric data should also resemble human data in this aspect.

**Discriminant Validity:** Discriminant validity ensures that distinct constructs are not excessively correlated. It can be examined by calculating the mean absolute correlation.

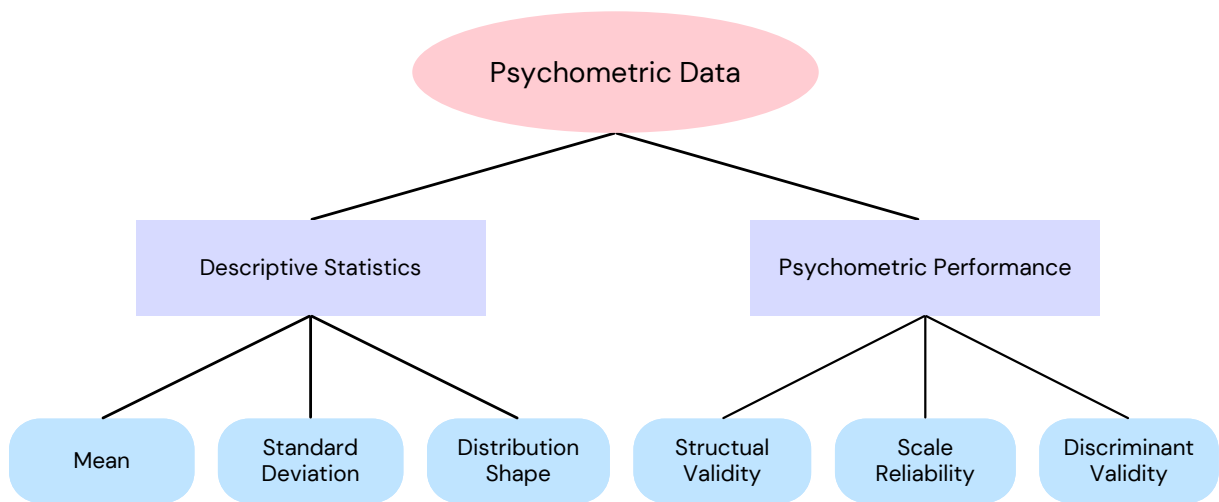


Figure 2: Psychometric Data Evaluation Framework



## I Appendix: Additional Analyses and Results

### I.1 Additional Method Comparison Results

**Additional Descriptive Statistics Results** Here we show the detailed means and standard deviations for human responses and LLM responses at the item, facet, and domain levels, see Tables 12, 13, 14, 15, 16, and 17.

#### **Additional Psychometric Performance Results**

Here we show the TFM fit information in Table 18; FFM fit information in Table 19; TCC results for the TFM of each BFI-2 domain in Table 20; TCC results for the FFM in Table 21; specific standardized factor loading results in Tables 22 and 23; specific inter-factor correlation results in Tables 24 and 25.

**Scale Reliability:** Facet level and domain level Cronbach's alpha for different method LLM responses on BFI-2 and human responses are shown in Table 26.

It can be observed that the PSI method, compared to the Persona and Shape methods, performs closer to the results of the human sample in terms of Cronbach's alpha (with the number of data marked in italics and bold being the smallest). This indicates that the LLM personality data generated by the PSI method holds an advantage in consistency and reliability, enabling it to more accurately simulate the statistical characteristics of human samples.

**Discriminant Validity:** The results for discriminant validity are shown in Table 27. We can further observe that, compared to the Persona and Shape methods, the PSI method demonstrates the closest performance to human samples in terms of discriminant validity. Specifically, the PSI method shows the closest mean of absolute values when examining its correlations with human samples.

Both higher and lower levels of external validity reveal the degree of differences between the methods and human samples. Higher external validity indicates that the Big Five factors in human samples are more distinctly differentiated from one another, while lower external validity suggests the opposite. Therefore, our focus here is on identifying the approach that most closely aligns with the performance of human samples.

No.	Item Content	Human	Persona		Shape		PSI	
			GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
1	I am someone who is outgoing, sociable	3.03	3.41	2.77	2.76	2.45	2.67	2.14
2	I am someone who is compassionate, has a soft heart	4.34	3.51	3.87	3.10	3.17	4.05	3.91
3	I am someone who tends to be disorganized	3.62	3.39	3.77	3.32	3.31	3.34	3.24
4	I am someone who is relaxed, handles stress well	2.56	3.20	2.94	3.47	3.32	2.91	2.96
5	I am someone who has few artistic interests	3.68	3.72	3.80	3.27	2.92	2.75	2.73
6	I am someone who has an assertive personality	2.96	3.08	3.07	2.82	2.81	2.57	2.25
7	I am someone who is respectful, treats others with respect	4.62	3.38	4.43	3.48	3.62	4.16	4.22
8	I am someone who tends to be lazy	3.62	4.00	4.23	3.51	3.44	3.54	3.75
9	I am someone who stays optimistic after a setback	2.44	2.74	2.29	2.95	2.92	2.18	2.29
10	I am someone who is curious about many different things	4.43	3.89	3.71	3.02	3.15	3.60	3.05
11	I am someone who rarely feels excited or eager	3.54	4.30	4.18	3.21	2.99	2.98	2.40
12	I am someone who tends to find fault with others	3.34	3.29	4.31	3.47	3.46	3.58	3.71
13	I am someone who is dependable, steady	4.33	3.29	3.70	2.93	3.11	4.12	3.96
14	I am someone who is moody, has up and down mood swings	2.56	3.01	2.52	2.72	2.60	3.04	2.31
15	I am someone who is inventive, finds clever ways to do things	3.90	3.44	3.24	3.06	2.91	3.05	2.37
16	I am someone who tends to be quiet	2.35	3.26	3.44	3.02	3.25	2.17	2.51
17	I am someone who feels little sympathy for others	3.71	3.73	4.47	3.72	3.54	4.22	2.88
18	I am someone who is systematic, likes to keep things in order	3.96	3.00	2.46	2.94	2.59	3.18	2.44
19	I am someone who can be tense	3.29	3.23	3.18	3.12	3.04	3.66	2.73
20	I am someone who is fascinated by art, music, or literature	4.06	3.61	2.99	3.15	2.90	3.43	2.81
21	I am someone who has an assertive personality (dominant)	2.90	2.94	2.46	2.77	2.53	2.42	2.07
22	I am someone who starts arguments with others	4.23	3.65	4.56	3.65	3.70	4.44	4.35
23	I am someone who has difficulty getting started on tasks	3.31	3.29	3.85	3.16	3.23	3.10	2.79
24	I am someone who feels secure, comfortable with self	2.26	2.75	2.24	3.14	2.65	2.65	2.70
25	I am someone who avoids intellectual, philosophical discussions	3.93	3.40	3.36	3.23	3.13	2.85	2.57
26	I am someone who is less active than other people	3.35	3.80	3.86	3.36	3.14	2.65	2.67
27	I am someone who has a forgiving nature	3.78	3.05	3.48	3.05	3.28	3.51	3.59
28	I am someone who can be somewhat careless	3.56	3.30	3.19	3.16	3.00	3.06	3.15
29	I am someone who is emotionally stable, not easily upset	2.44	3.16	2.66	3.17	2.87	2.70	2.67
30	I am someone who has little creativity	3.95	4.10	4.36	3.59	3.42	3.37	2.85
31	I am someone who is sometimes shy, introverted	2.35	3.28	3.43	3.11	3.34	2.14	2.70
32	I am someone who is helpful and unselfish with others	4.18	3.32	3.88	3.27	3.24	4.06	4.12
33	I am someone who keeps things neat and tidy	3.72	2.99	2.47	3.07	2.63	3.04	2.70
34	I am someone who worries a lot	3.30	2.96	2.59	2.79	2.50	3.64	3.04
35	I am someone who values art and beauty	4.16	3.65	3.23	3.14	3.06	3.38	2.87
36	I am someone who finds it hard to influence people	3.09	3.32	3.44	3.25	2.99	2.74	2.25
37	I am someone who is sometimes rude to others	3.79	3.24	4.43	3.50	3.63	3.70	4.10
38	I am someone who is efficient, gets things done	4.23	3.31	3.45	3.07	3.12	3.69	3.07
39	I am someone who often feels sad	2.71	2.86	2.27	2.74	2.42	3.24	2.69
40	I am someone who is complex, a deep thinker	4.02	3.64	3.59	3.19	3.14	3.34	2.74
41	I am someone who is full of energy	3.14	3.45	3.39	2.99	2.86	2.81	2.41
42	I am someone who is suspicious of others' intentions	2.75	3.14	3.91	3.27	3.40	3.42	3.60
43	I am someone who is reliable, can always be counted on	4.32	3.28	3.76	3.12	3.18	4.06	3.75
44	I am someone who keeps their emotions under control	2.22	3.10	2.92	2.98	2.90	2.50	2.34
45	I am someone who has difficulty imagining things	4.17	4.06	4.08	3.57	3.43	3.21	2.14
46	I am someone who is talkative	2.88	3.11	2.53	2.87	2.50	3.31	2.36
47	I am someone who can be cold and uncaring	3.84	3.95	4.44	3.60	3.47	4.08	4.19
48	I am someone who leaves a mess, doesn't clean up	4.15	3.47	4.25	3.39	3.25	3.32	4.03
49	I am someone who rarely feels anxious or afraid	3.44	3.20	3.12	2.85	2.71	3.48	3.66
50	I am someone who thinks poetry and plays are boring	3.66	3.54	3.85	3.44	3.15	2.90	2.29
51	I am someone who prefers to have others take charge	2.96	3.54	4.33	3.39	3.20	3.52	3.87
52	I am someone who is polite, courteous to others	4.50	3.28	4.13	3.33	3.56	3.97	4.38
53	I am someone who is persistent, works until the task is finished	4.27	3.56	3.87	3.01	3.13	4.07	3.75
54	I am someone who tends to feel depressed, blue	2.65	2.70	2.01	2.60	2.36	3.25	2.78
55	I am someone who has little interest in abstract ideas	3.83	3.25	3.40	3.08	2.93	2.53	2.15
56	I am someone who shows a lot of enthusiasm	3.35	3.86	4.08	2.81	2.81	2.70	2.58
57	I am someone who assumes the best about people	3.37	3.12	3.40	3.04	3.13	3.46	3.64
58	I am someone who sometimes behaves irresponsibly	3.64	3.11	3.44	2.90	3.22	3.13	3.55
59	I am someone who is temperamental, gets emotional easily	2.53	3.02	2.33	2.80	2.66	2.96	2.04
60	I am someone who is original, comes up with new ideas	3.84	3.28	3.25	2.96	2.88	3.05	2.40

Table 12: Item Level Mean for BFI-2 Human Responses and Different Methods LLM Responses

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

No.	Item Content	Human	Persona		Shape		PSI	
			GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
1	I am someone who is outgoing, sociable	1.40	0.87	1.31	1.22	1.37	0.98	0.78
2	I am someone who is compassionate, has a soft heart	0.90	0.84	1.12	1.22	1.44	0.93	1.03
3	I am someone who tends to be disorganized	1.36	0.78	0.95	1.17	1.46	0.76	0.96
4	I am someone who is relaxed, handles stress well	1.30	0.79	1.16	1.16	1.39	0.92	0.96
5	I am someone who has few artistic interests	1.33	1.04	1.31	1.25	1.32	1.39	1.11
6	I am someone who has an assertive personality	1.38	0.79	1.18	1.18	1.37	0.76	0.73
7	I am someone who is respectful, treats others with respect	0.61	0.89	0.92	1.23	1.42	0.75	0.78
8	I am someone who tends to be lazy	1.31	0.87	1.04	1.28	1.50	0.81	0.97
9	I am someone who stays optimistic after experiencing a setback	1.24	0.84	0.98	1.34	1.54	0.89	0.93
10	I am someone who is curious about many different things	0.81	0.89	1.12	1.37	1.55	0.98	0.91
11	I am someone who rarely feels excited or eager	1.26	0.84	0.75	1.33	1.41	0.93	0.82
12	I am someone who tends to find fault with others	1.32	0.59	0.69	1.26	1.48	0.83	0.70
13	I am someone who is dependable, steady	0.86	0.85	1.12	1.33	1.59	0.79	0.75
14	I am someone who is moody, has up and down mood swings	1.36	0.56	1.09	1.27	1.41	0.76	0.93
15	I am someone who is inventive, finds clever ways to do things	1.03	0.69	1.14	1.21	1.44	0.63	0.74
16	I am someone who tends to be quiet	1.34	0.84	1.12	1.12	1.27	1.02	0.97
17	I am someone who feels little sympathy for others	1.45	0.93	0.79	1.27	1.54	0.68	1.00
18	I am someone who is systematic, likes to keep things in order	1.06	0.88	1.09	1.16	1.37	0.94	0.81
19	I am someone who can be tense	1.27	0.75	1.08	1.22	1.20	0.95	0.95
20	I am someone who is fascinated by art, music, or literature	1.17	0.89	1.44	1.18	1.36	1.18	1.11
21	I am someone who is dominant, acts as a leader	1.35	0.74	1.15	1.13	1.49	0.70	0.61
22	I am someone who starts arguments with others	1.05	0.72	0.77	1.24	1.54	0.72	0.60
23	I am someone who has difficulty getting started on tasks	1.39	0.69	0.86	1.26	1.49	0.74	0.69
24	I am someone who feels secure, comfortable with self	1.27	0.88	1.08	1.42	1.49	0.96	0.92
25	I am someone who avoids intellectual, philosophical discussions	1.21	0.79	1.36	1.25	1.56	1.10	0.99
26	I am someone who is less active than other people	1.31	0.99	1.04	1.24	1.31	0.98	0.83
27	I am someone who has a forgiving nature	1.20	0.57	0.90	1.20	1.37	0.85	0.78
28	I am someone who can be somewhat careless	1.25	0.83	1.11	1.18	1.38	0.68	1.00
29	I am someone who is emotionally stable, not easily upset	1.27	0.70	1.01	1.31	1.47	1.06	1.09
30	I am someone who has little creativity	1.20	0.76	0.69	1.22	1.32	0.83	1.02
31	I am someone who is sometimes shy, introverted	1.36	0.87	1.11	1.07	1.30	1.04	0.99
32	I am someone who is helpful and unselfish with others	0.88	0.80	1.05	1.21	1.50	0.73	0.72
33	I am someone who keeps things neat and tidy	1.24	0.64	1.00	1.09	1.35	0.73	0.71
34	I am someone who worries a lot	1.47	0.70	1.08	1.23	1.33	1.18	1.07
35	I am someone who values art and beauty	1.07	0.82	1.30	1.18	1.35	1.00	1.06
36	I am someone who finds it hard to influence people	1.20	0.69	1.04	1.26	1.37	0.67	0.67
37	I am someone who is sometimes rude to others	1.25	0.66	0.91	1.31	1.49	0.77	1.00
38	I am someone who is efficient, gets things done	0.88	0.82	1.12	1.25	1.46	0.84	0.91
39	I am someone who often feels sad	1.43	0.68	0.98	1.27	1.25	0.99	0.91
40	I am someone who is complex, a deep thinker	1.08	0.72	1.06	1.06	1.26	0.83	0.99
41	I am someone who is full of energy	1.30	0.89	1.24	1.23	1.35	0.85	0.76
42	I am someone who is suspicious of others' intentions	1.30	0.58	0.83	1.26	1.49	1.00	0.68
43	I am someone who is reliable, can always be counted on	0.88	0.93	1.09	1.32	1.59	0.80	0.78
44	I am someone who keeps their emotions under control	1.14	0.64	1.08	1.22	1.44	0.77	1.04
45	I am someone who has difficulty imagining things	1.06	0.77	0.53	1.25	1.27	0.95	0.54
46	I am someone who is talkative	1.38	0.70	1.11	1.04	1.28	1.09	0.92
47	I am someone who can be cold and uncaring	1.23	0.80	0.87	1.35	1.49	0.78	0.98
48	I am someone who leaves a mess, doesn't clean up	1.12	0.71	0.95	1.26	1.64	0.55	0.85
49	I am someone who rarely feels anxious or afraid	1.37	0.83	1.12	1.15	1.37	1.02	0.89
50	I am someone who thinks poetry and plays are boring	1.37	0.75	1.31	1.38	1.53	0.87	0.68
51	I am someone who prefers to have others take charge	1.28	0.69	0.76	1.27	1.57	0.70	0.51
52	I am someone who is polite, courteous to others	0.70	0.65	0.97	1.14	1.44	0.75	0.79
53	I am someone who is persistent, works until the task is finished	0.91	0.87	1.13	1.26	1.48	0.80	0.94
54	I am someone who tends to feel depressed, blue	1.46	0.79	1.07	1.31	1.37	1.08	1.00
55	I am someone who has little interest in abstract ideas	1.22	0.96	1.08	1.22	1.36	1.17	0.63
56	I am someone who shows a lot of enthusiasm	1.22	0.85	1.09	1.34	1.51	0.83	0.94
57	I am someone who assumes the best about people	1.28	0.78	0.98	1.27	1.47	0.92	0.74
58	I am someone who sometimes behaves irresponsibly	1.30	0.74	1.13	1.14	1.40	0.70	0.95
59	I am someone who is temperamental, gets emotional easily	1.33	0.64	1.08	1.24	1.52	1.00	0.87
60	I am someone who is original, comes up with new ideas	1.07	0.65	1.13	1.22	1.41	0.68	0.75

Table 13: Item Level Standard Deviation for BFI-2 Human Responses and Different Methods LLM Responses  
*Note:*  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Facet	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Sociability	2.65	3.26	3.04	2.94	2.89	2.57	2.43
Assertiveness	2.97	3.22	3.33	3.06	2.88	2.81	2.61
Energy Level	3.35	3.85	3.88	3.09	2.95	2.79	2.52
Compassion	4.02	3.63	4.16	3.42	3.36	4.10	3.78
Respectfulness	4.29	3.39	4.39	3.49	3.63	4.07	4.26
Trust	3.31	3.15	3.78	3.21	3.32	3.49	3.64
Organization	3.86	3.21	3.24	3.18	2.94	3.22	3.10
Productiveness	3.86	3.54	3.85	3.19	3.23	3.60	3.34
Responsibility	3.96	3.25	3.52	3.03	3.13	3.59	3.60
Anxiety	3.15	3.15	2.96	3.06	2.89	3.42	3.10
Depression	2.51	2.76	2.20	2.86	2.59	2.83	2.61
Emotional Volatility	2.44	3.07	2.61	2.92	2.76	2.80	2.34
Intellectual Curiosity	4.05	3.55	3.51	3.13	3.09	3.08	2.63
Aesthetic Sensitivity	3.89	3.63	3.47	3.25	3.01	3.12	2.67
Creative Imagination	3.96	3.72	3.73	3.29	3.16	3.17	2.44

Table 14: Facet Level Mean for BFI-2 Human Responses and Different Methods LLM Responses

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Facet	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Sociability	1.17	0.69	0.98	0.95	1.00	0.91	0.74
Assertiveness	1.04	0.55	0.82	1.10	1.22	0.57	0.42
Energy Level	0.95	0.75	0.88	1.12	1.19	0.77	0.64
Compassion	0.81	0.69	0.81	1.18	1.36	0.69	0.76
Respectfulness	0.70	0.58	0.74	1.13	1.33	0.61	0.68
Trust	1.00	0.44	0.67	1.17	1.34	0.77	0.62
Organization	1.01	0.57	0.79	0.98	1.18	0.63	0.68
Productiveness	0.91	0.65	0.87	1.12	1.29	0.68	0.68
Responsibility	0.84	0.59	0.95	1.10	1.26	0.64	0.74
Anxiety	1.11	0.54	0.89	0.95	1.04	0.88	0.75
Depression	1.14	0.56	0.89	1.16	1.26	0.87	0.81
Emotional Volatility	1.10	0.43	0.92	1.06	1.25	0.76	0.87
Intellectual Curiosity	0.83	0.68	0.96	1.04	1.24	0.89	0.70
Aesthetic Sensitivity	1.01	0.75	1.09	1.12	1.26	0.96	0.83
Creative Imagination	0.88	0.57	0.74	1.05	1.17	0.67	0.60

Table 15: Facet Level Standard Deviation for BFI-2 Human Responses and Different Methods LLM Responses

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Domain	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Extraversion	2.99	3.44	3.42	3.03	2.90	2.72	2.52
Agreeableness	3.87	3.39	4.11	3.37	3.43	3.89	3.89
Conscientiousness	3.89	3.33	3.54	3.13	3.10	3.47	3.35
Neuroticism	2.70	2.99	2.59	2.94	2.75	3.02	2.68
Openness	3.97	3.63	3.57	3.22	3.09	3.12	2.58

Table 16: Domain Level Mean for BFI-2 Human Responses and Different Methods LLM Responses  
*Note:*  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Domain	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Extraversion	0.85	0.57	0.81	0.95	1.00	0.67	0.52
Agreeableness	0.69	0.51	0.70	1.11	1.28	0.62	0.62
Conscientiousness	0.80	0.53	0.79	1.00	1.15	0.59	0.64
Neuroticism	1.02	0.43	0.81	0.97	1.09	0.77	0.72
Openness	0.76	0.58	0.78	1.03	1.15	0.72	0.59

Table 17: Domain Level Standard Deviation for BFI-2 Human Responses and Different Methods LLM Responses  
*Note:*  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.



Domain	Model	$\chi^2$	<i>df</i>	CFI	TLI	RMSEA	SRMR
Extraversion	Human	993.189	51	.892	.861	.109	.057
				Persona			
	GPT-4o	261.099	51	.886	.852	.117	.058
	Llama3	316.715	51	.885	.851	.132	.060
				Shape			
	GPT-4o	663.947	51	.823	.771	.200	.101
	Llama3	934.002	51	.684	.592	.241	.132
				PSI			
	GPT-4o	329.341	51	.892	.860	.135	.057
	Llama3	332.944	51	.810	.755	.136	.092
Agreeableness	Human	904.640	51	.875	.838	.104	.060
				Persona			
	GPT-4o	210.863	51	.900	.870	.102	.062
	Llama3	616.199	51	.798	.739	.192	.088
				Shape			
	GPT-4o	525.799	51	.909	.882	.176	.039
	Llama3	874.914	51	.840	.793	.233	.070
				PSI			
	GPT-4o	177.494	51	.952	.937	.091	.049
	Llama3	255.684	51	.923	.900	.116	.050
Conscientiousness	Human	1041.784	51	.897	.867	.112	.058
				Persona			
	GPT-4o	385.339	51	.780	.716	.148	.088
	Llama3	705.359	51	.768	.699	.207	.109
				Shape			
	GPT-4o	705.340	51	.826	.775	.207	.069
	Llama3	1502.901	51	.642	.536	.310	.196
				PSI			
	GPT-4o	348.206	51	.887	.854	.139	.052
	Llama3	427.908	51	.834	.785	.157	.069
Neuroticism	Human	929.871	51	.931	.911	.105	.052
				Persona			
	GPT-4o	304.611	51	.756	.684	.129	.091
	Llama3	346.207	51	.884	.850	.139	.059
				Shape			
	GPT-4o	960.437	51	.720	.638	.244	.099
	Llama3	808.305	51	.772	.705	.224	.101
				PSI			
	GPT-4o	448.933	51	.876	.840	.161	.069
	Llama3	448.871	51	.854	.811	.161	.067
Openness	Human	909.210	51	.899	.870	.104	.064
				Persona			
	GPT-4o	226.649	51	.910	.883	.107	.057
	Llama3	256.032	51	.904	.876	.116	.077
				Shape			
	GPT-4o	502.104	51	.878	.842	.172	.055
	Llama3	569.466	51	.861	.820	.185	.057
				PSI			
	GPT-4o	310.479	51	.904	.876	.130	.071
	Llama3	352.224	51	.854	.811	.140	.114

Table 18: Model Fits for BFI-2 Three-Factor Models of Each Domain

*Note:*  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Model	$\chi^2$	df	CFI	TLI	RMSEA	SRMR
Human	1747.433	80	.852	.806	.116	.080
Persona						
GPT-4o	551.459	80	.819	.762	.140	.097
Llama3	954.848	80	.780	.711	.191	.120
Shape						
GPT-4o	2132.208	80	.711	.621	.292	.143
Llama3	1988.584	80	.717	.629	.283	.149
PSI						
GPT-4o	899.057	80	.769	.697	.185	.119
Llama3	802.352	80	.774	.703	.173	.114

Table 19: Model Fits for BFI-2 Five-Factor Model

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Domain	Facet	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Extraversion	Sociability	1.00	.99	1.00	.99	1.00	.99
	Assertiveness	.97	.98	.99	.97	.98	.97
	Energy Level	.98	.98	.97	.96	.98	.98
Agreeableness	Compassion	.95	.96	.97	.95	.97	1.00
	Respectfulness	1.00	1.00	1.00	1.00	.98	.99
	Trust	.98	.99	1.00	.99	1.00	.99
Conscientiousness	Organization	.99	1.00	.99	.95	.99	1.00
	Productiveness	.99	.99	.99	.99	1.00	.97
	Responsibility	.99	.98	.99	.96	1.00	.99
Neuroticism	Anxiety	.99	.99	.98	.99	1.00	1.00
	Depression	.97	.99	.99	.99	1.00	1.00
	Emotional Volatility	.97	1.00	.99	.99	1.00	1.00
Openness	Intellectual Curiosity	1.00	1.00	.99	.98	1.00	.98
	Aesthetic Sensitivity	.99	.98	.99	.98	.99	.97
	Creative Imagination	1.00	.98	1.00	1.00	.99	.96

Table 20: Tucker's Congruence Coefficient for BFI-2 Three-Factor Models of Each Domain

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Domain	Persona		Shape		PSI	
	GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Extraversion	1.00	.99	1.00	1.00	1.00	1.00
Agreeableness	1.00	1.00	.99	.99	1.00	1.00
Conscientiousness	.99	.98	.99	.99	.99	.99
Neuroticism	.99	1.00	.99	1.00	1.00	1.00
Openness	1.00	.99	1.00	.99	1.00	.99

Table 21: Tucker's Congruence Coefficient for BFI-2 Five-Factor Model

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

Domain	Facet--Item	Human	Persona		Shape		PSI	
			GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Extraversion	Sociability--item1	.80	.77	.72	.85	<b>.53</b>	.87	.63
	Sociability--item16	.82	.79	.88	.74	.75	.84	.88
	Sociability--item31	.82	.79	.87	.73	.80	.91	.76
	Sociability--item46	.76	.74	.64	.85	.57	.76	.67
	Assertiveness--item6	.76	.66	.84	.94	.85	.81	.69
	Assertiveness--item21	.89	<b>.61</b>	.73	.90	<b>.69</b>	.80	<b>.69</b>
	Assertiveness--item36	.53	<b>.73</b>	.72	<b>.84</b>	<b>.88</b>	<b>.76</b>	.44
	Assertiveness--item51	.69	.59	.54	.80	.73	.55	<b>.27</b>
	Energy Level--item11	.45	<b>.82</b>	<b>.73</b>	<b>.91</b>	<b>.89</b>	<b>.82</b>	<b>.66</b>
	Energy Level--item26	.54	.69	<b>.75</b>	.72	<b>.78</b>	<b>.74</b>	.47
	Energy Level--item41	.83	.81	.87	.82	.76	.87	.71
Energy Level--item56	.77	.82	.78	.88	.77	.80	.84	
Agreeableness	Compassion--item2	.73	.76	.75	.90	.74	.84	.83
	Compassion--item17	.34	<b>.84</b>	<b>.86</b>	<b>.93</b>	<b>.97</b>	<b>.79</b>	.46
	Compassion--item32	.62	.74	.74	<b>.87</b>	<b>.86</b>	.80	.79
	Compassion--item47	.75	.67	.85	<b>.96</b>	<b>.98</b>	.89	.86
	Respectfulness--item7	.71	.82	.80	<b>.93</b>	.84	.90	<b>.91</b>
	Respectfulness--item22	.60	.54	.68	<b>.80</b>	<b>.87</b>	.43	.48
	Respectfulness--item37	.70	.74	.79	<b>.94</b>	<b>.94</b>	.80	.82
	Respectfulness--item52	.70	.70	.80	.83	.81	.85	.86
	Trust--item12	.69	.63	<b>.35</b>	.57	<b>.95</b>	<b>.95</b>	.78
	Trust--item27	.68	.60	.80	.84	.84	.79	.84
Trust--item42	.67	<b>.35</b>	.63	<b>.95</b>	<b>.95</b>	.78	.83	
Trust--item57	.78	.64	.82	.90	.85	.88	.89	
Conscientiousness	Organization--item3	.86	<b>.64</b>	.75	.75	<b>.48</b>	.83	.83
	Organization--item18	.70	.62	.66	.73	<b>.90</b>	.83	.65
	Organization--item33	.86	<b>.63</b>	.68	.95	<b>.55</b>	.74	.74
	Organization--item48	.74	.72	.72	.68	.55	.80	.79
	Productiveness--item8	.66	.67	.68	.64	.67	.73	<b>.43</b>
	Productiveness--item23	.86	.71	<b>.65</b>	.81	.70	.78	<b>.43</b>
	Productiveness--item38	.87	.72	.88	.82	.85	.85	<b>.52</b>
	Productiveness--item53	.65	.47	.66	.80	.63	.81	.82
	Responsibility--item13	.72	.67	<b>.94</b>	<b>.92</b>	<b>.97</b>	.79	.83
	Responsibility--item28	.77	.58	.58	.68	<b>.46</b>	.59	.71
	Responsibility--item43	.66	.51	.74	.75	.70	.68	.57
Responsibility--item58	.80	<b>.59</b>	.80	.75	.87	.71	.74	
Neuroticism	Anxiety--item4	.60	.71	.72	.68	.62	<b>.80</b>	.79
	Anxiety--item19	.63	.65	.61	.65	.72	.78	<b>.43</b>
	Anxiety--item34	.65	.47	.81	.73	.79	.59	.65
	Anxiety--item49	.74	.71	.84	.85	.70	.92	.89
	Depression--item9	.66	.67	.60	.81	.77	<b>.88</b>	<b>.44</b>
	Depression--item24	.86	.82	<b>.59</b>	.81	.78	.92	.91
	Depression--item39	.87	.86	<b>.65</b>	<b>.47</b>	.86	.95	.95
	Depression--item54	.65	.69	<b>.88</b>	<b>.92</b>	<b>.95</b>	.79	<b>.85</b>
	Emotional Volatility--item14	.72	.77	.72	.59	.88	.84	.86
	Emotional Volatility--item29	.77	.66	<b>.51</b>	.74	.75	.70	.68
Emotional Volatility--item44	.66	.64	<b>.27</b>	.68	.75	.60	.67	
Emotional Volatility--item59	.80	.67	.72	.83	.79	.65	.85	
Openness	Intellectual Curiosity--item10	.60	<b>.90</b>	<b>.81</b>	<b>.84</b>	<b>.91</b>	<b>.92</b>	<b>.89</b>
	Intellectual Curiosity--item25	.63	<b>.91</b>	.67	<b>.88</b>	<b>.89</b>	<b>.95</b>	<b>.94</b>
	Intellectual Curiosity--item40	.65	.80	.59	.70	.57	<b>.85</b>	<b>.85</b>
	Intellectual Curiosity--item55	.74	.77	.81	.84	.89	.88	.56
	Aesthetic Sensitivity--item5	.66	.66	.84	.85	<b>.91</b>	.71	.74
	Aesthetic Sensitivity--item20	.86	.77	.81	<b>.65</b>	.83	.79	.85
	Aesthetic Sensitivity--item35	.87	<b>.59</b>	.70	.84	.81	.92	.88
	Aesthetic Sensitivity--item50	.65	.63	.69	<b>.41</b>	<b>.91</b>	.63	<b>.32</b>
	Creative Imagination--item15	.72	.77	<b>.40</b>	.71	.81	<b>.95</b>	.79
	Creative Imagination--item30	.77	.80	.59	.70	.81	.86	.82
Creative Imagination--item45	.66	.84	.68	.49	.77	<b>.92</b>	<b>.35</b>	
Creative Imagination--item60	.80	.86	.70	.95	.81	.84	.95	

Table 22: Standardized Factor Loadings for BFI-2 Three-Factor Models of Each Domain with Human Responses  
*Note:*  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis. Italics for absolute differences compared to the human responses of .100 to .199, and boldface for differences of .200 or higher.

Domain--Facet	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
E--Sociability	.66	.76	.85	<b>.91</b>	<b>.87</b>	.81	.78
E--Assertiveness	.59	.67	<b>.86</b>	.77	.73	.71	.68
E--Energy Level	.77	.88	.86	.89	.85	.92	.89
A--Compassion	.70	.82	<b>.90</b>	<b>.96</b>	<b>.99</b>	.82	.79
A--Respectfulness	.80	.90	.94	.90	.85	.88	.94
A--Trust	.65	.75	<b>.87</b>	<b>.93</b>	<b>.95</b>	.84	.84
C--Organization	.70	.87	.87	<b>.94</b>	<b>.97</b>	.87	.86
C--Productiveness	.89	.77	.75	.84	.81	.84	.82
C--Responsibility	.79	.82	.95	.93	.92	.88	.91
N--Anxiety	.84	.69	.75	.80	.86	.77	.77
N--Depression	.88	.91	.96	1.02	.96	.97	.87
N--Emotional Volatility	.85	<b>.63</b>	.79	.71	.80	.78	.85
O--Intellectual Curiosity	.74	.77	.66	<b>.97</b>	.93	.74	.85
O--Aesthetic Sensitivity	.71	.73	.65	<b>.95</b>	<b>.96</b>	.69	.59
O--Creative Imagination	.81	.91	.89	.90	.82	.96	.76

Table 23: Standardized Factor Loadings for BFI-2 Five-Factor Model

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis. Italics for absolute differences compared to the human responses of .100 to .199, and boldface for differences of .200 or higher.

Domain	Facet--Facet	Human	Persona		Shape		PSI	
			GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Extraversion	Sociability--Assertiveness	.59	.72	<b>.89</b>	.77	<b>.91</b>	.69	.70
	Sociability--Energy Level	.64	.79	<b>.84</b>	<b>.91</b>	.80	<b>.84</b>	<b>.85</b>
	Assertiveness--Energy Level	.47	.66	<b>.82</b>	.65	.61	<b>.72</b>	.65
Agreeableness	Compassion--Respectfulness	.81	.90	.98	.92	.87	.85	.90
	Compassion--Trust	.70	.83	<b>.94</b>	<b>.94</b>	<b>.96</b>	.79	.84
	Respectfulness--Trust	.58	<b>.95</b>	<b>.96</b>	<b>.93</b>	<b>.84</b>	<b>.85</b>	<b>.86</b>
Conscientiousness	Organization--Productiveness	.75	.86	.81	<b>.97</b>	.70	.86	.86
	Organization--Responsibility	.70	<b>1.01</b>	.85	<b>.99</b>	.86	.89	.87
	Productiveness--Responsibility	.86	.82	.87	.87	.89	.82	.87
Neuroticism	Anxiety--Depression	.81	.88	.88	.95	.96	.86	.79
	Anxiety--Emotional Volatility	.90	.87	.79	.87	.87	.92	.83
	Depression--Emotional Volatility	.79	.84	.85	.89	.83	.86	.82
Openness	Intellectual Curiosity--Aesthetic Sensitivity	.66	.68	.59	<b>1.02</b>	<b>1.00</b>	.51	.51
	Intellectual Curiosity--Creative Imagination	.74	.85	.67	<b>.99</b>	.89	.78	.72
	Aesthetic Sensitivity--Creative Imagination	.63	.79	.60	<b>.93</b>	<b>.86</b>	.77	.50

Table 24: Inter-factor Correlations for BFI-2 Three-Factor Models of Each Domain with Human Responses

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis. Italics for absolute differences compared to the human responses of .100 to .199, and boldface for differences of .200 or higher.

Domain~~Domain	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Extraversion~~Agreeableness	.28	.26	.10	<b>.59</b>	<b>.54</b>	.41	.24
Extraversion~~Conscientiousness	.56	.43	<b>.26</b>	<b>.36</b>	<b>.31</b>	.44	<b>.34</b>
Extraversion~~Neuroticism	-.61	-.61	-.59	-.80	-.72	-.62	<b>-.31</b>
Extraversion~~Openness	.35	.37	.51	<b>.72</b>	<b>.71</b>	.29	.48
Agreeableness~~Conscientiousness	.47	.56	<b>.77</b>	.65	.65	.59	.63
Agreeableness~~Neuroticism	-.43	-.40	-.54	<b>-.66</b>	<b>-.73</b>	-.44	-.47
Agreeableness~~Openness	.27	.31	.25	<b>.79</b>	<b>.89</b>	.14	.14
Conscientiousness~~Neuroticism	-.61	-.52	-.63	-.67	<b>-.82</b>	-.58	-.60
Conscientiousness~~Openness	.17	.10	.16	.27	<b>.45</b>	.06	.18
Neuroticism~~Openness	-.17	-.06	-.23	<b>-.48</b>	<b>-.65</b>	<b>.14</b>	<b>.18</b>

Table 25: Inter-factor Correlations for BFI-2 Five-Factor Models with Human Responses

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis. Italics for absolute differences compared to the human responses of .100 to .199, and boldface for differences of .200 or higher.

Facet / Domain	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
Sociability	.87	.85	.86	.87	.77	.91	.82
Assertiveness	.81	.74	.80	.93	.86	.82	<b>.57</b>
Energy Level	.73	.86	.86	.90	.87	.88	.76
Compassion	.67	.84	<b>.87</b>	<b>.95</b>	<b>.93</b>	<b>.90</b>	.82
Respectfulness	.74	.79	.85	<b>.94</b>	.93	.84	.86
Trust	.80	.64	.79	.96	.94	.87	.88
Organization	.87	.73	.80	.86	.82	.86	.82
Productiveness	.80	.80	.85	.91	.90	.88	.76
Responsibility	.77	.67	.88	.91	.86	.88	.86
Anxiety	.84	.65	.81	.81	.79	.88	.78
Depression	.86	<b>.66</b>	.88	.90	.91	.91	.88
Emotional Volatility	.89	<b>.60</b>	.89	.87	.88	.89	.91
Intellectual Curiosity	.75	.82	.85	.87	.89	.89	.79
Aesthetic Sensitivity	.83	.87	.83	.92	.93	.87	.84
Creative Imagination	.82	.80	.82	.88	.89	.88	.77
Extraversion	.87	.90	.93	.95	.92	.93	.87
Agreeableness	.85	.90	.94	.98	.97	.93	.93
Conscientiousness	.90	.88	.93	.95	.94	.94	.92
Neuroticism	.94	.83	.93	.94	.94	.95	.93
Openness	.89	.91	.89	.96	.96	.92	.88

Table 26: Cronbach's alpha for BFI-2 Human Responses and Different Methods LLM Responses

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis. Italics for absolute differences compared to the human responses of .100 to .199, and boldface for differences of .200 or higher.



Dimensions	Human	Persona		Shape		PSI	
		GPT-4o	Llama3	GPT-4o	Llama3	GPT-4o	Llama3
E, A	.17	.21	.12	.48	.40	.33	.21
E, C	.35	.40	.33	.46	.41	.38	.30
E, N	-.45	-.43	-.42	-.58	-.57	-.43	-.21
E, O	.22	.23	.28	.66	.68	.13	.33
A, C	.33	.47	.66	.67	.70	.49	.52
A, N	-.36	-.30	-.54	-.64	-.71	-.29	-.37
A, O	.20	.30	.31	.71	.78	.14	.14
C, N	-.49	-.43	-.60	-.74	-.87	-.47	-.52
C, O	.09	.09	.20	.29	.48	.02	.08
N, O	-.11	.03	-.08	-.31	-.57	.32	.23
Mean of Absolute Values	.28	.29	.35	.55	.62	.30	.29

Table 27: Domain Level Correlation Analysis for BFI-2 Human Responses and Different Methods LLM Responses  
*Note:*  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis. E = Extraversion; A = Agreeableness; C = Conscientiousness; N = Neuroticism; O = Openness.

## I.2 Additional Personality-Related Behavioral Performance Results

Table 28 presents the correlations between the personality dimensions and OCB/CWB reported by PSI GPT-4o, with a comparison to human self-reported data. The relevant results for PSI Llama3 are provided in Table 29, while Figures 3 and Figure 4 display the complete correlation matrices for GPT-4o and Llama3, respectively.

Domain	OCB		CWB	
	Human	PSI	Human	PSI
Ext	.42	.54	-.01	-.08
Agr	.18	.36	-.30	-.54
Con	.12	.45	-.35	-.47
Neu	-.21	-.35	.23	.36
Ope	.17	.02	-.17	.01

Table 28: Comparison of OCB and CWB Correlations with Personality Domains: Human vs. PSI GPT-4o  
*Note: n = 357 for human and PSI method.*

Domain	OCB		CWB	
	Human	PSI	Human	PSI
Ext	.42	.48	-.01	.05
Agr	.18	.46	-.30	-.45
Con	.12	.47	-.35	-.41
Neu	-.21	-.38	.23	.33
Ope	.17	.11	-.17	.00

Table 29: Comparison of OCB and CWB Correlations with Personality Domains: Human vs. PSI Llama3  
*Note: n = 357 for human and PSI method.*

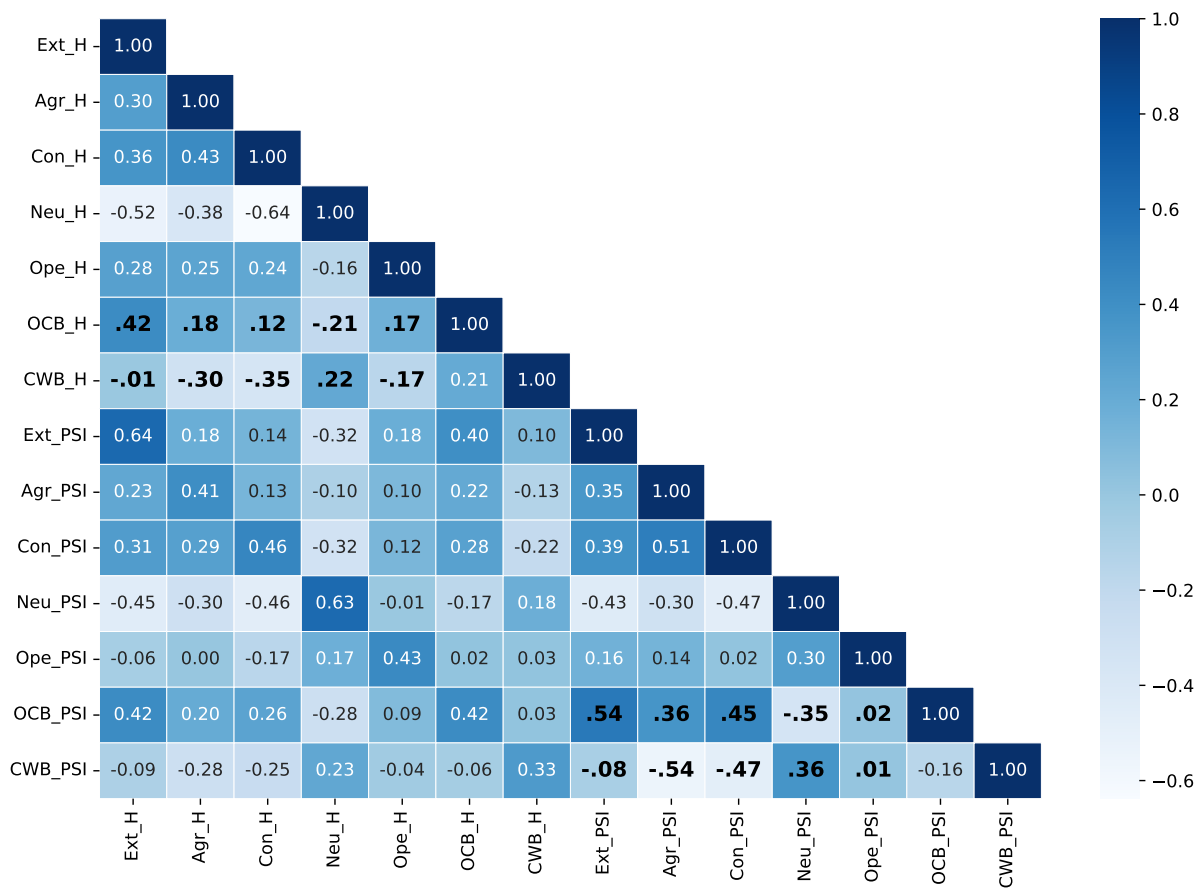


Figure 3: Complete Correlation Matrices for Human and PSI GPT-4o  
 Note:  $n = 357$  for human and PSI method.

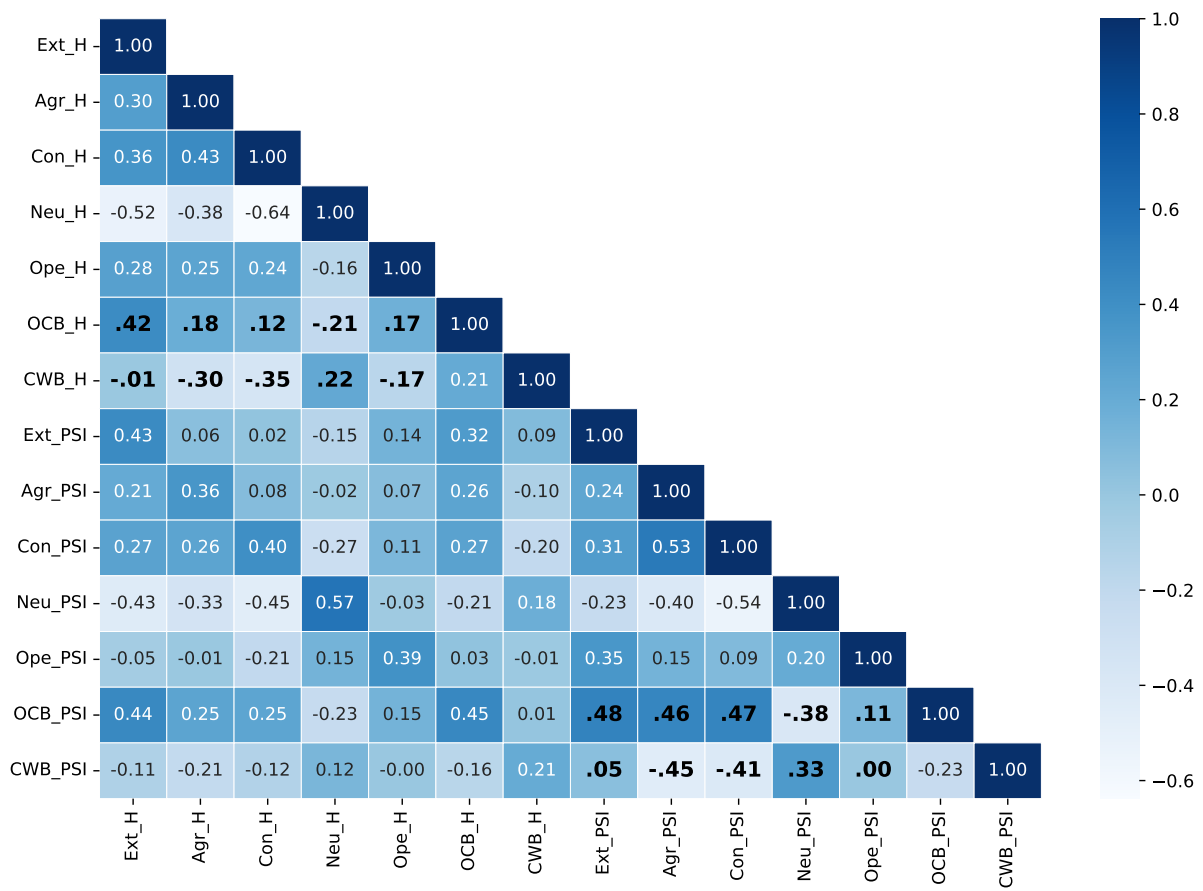


Figure 4: Complete Correlation Matrices for Human and PSI Llama3  
 Note:  $n = 357$  for human and PSI method.

### I.3 The Influence of Social Desirability

In previous studies, researchers have observed some intriguing phenomena, specifically that LLMs tend to exhibit human-like social desirability bias when simulating human samples. This bias is reflected in the LLMs' responses, which lean toward behaviors and traits that are socially approved.

Numerous studies support this finding. For example, Hilliard et al. (2024) noted that newer and larger-parameter LLMs display more diverse personality traits, including higher levels of agreeableness, emotional stability, and openness. Similarly, Salecha et al. (2024) found that LLMs exhibit human-like social desirability bias when generating simulated data.

In this section, using the data we have collected, we directly compare the results from human samples and LLM-simulated data. We examined the self-reported personality traits of human samples alongside the simulated personality traits from LLMs, analyzing how both were influenced by social desirability ratings. Specifically, we used social desirability ratings to predict and test the correlations and differences between the mean scores of human and LLM data across various items.

By analyzing the predictive relationships and correlations between item means and social desirability ratings, we aim to uncover how social desirability influences the characteristics and patterns in both human and LLM data. Furthermore, by comparing the differences in item means between human data and LLM-simulated data, we explore whether social desirability ratings affect human samples and model-simulated data in similar ways.

**Social Desirability Rating** Social desirability ratings of the BFI-2 items were obtained from another ongoing study where 142 human resource practitioners were asked to rate how desirable each item was in general (1 = "Very undesirable," 2 = "Undesirable," 3 = "Slightly undesirable," 4 = "Neither desirable nor undesirable," 5 = "Slightly desirable," 6 = "Desirable," 7 = "Very desirable").

If a certain measurement item receives a high score (6 or 7), it indicates that the trait is generally viewed as desirable or positive in society. On the other hand, a low score (1 or 2) suggests that the trait may be considered less ideal or inconsistent with social expectations. Some desirable traits may lead participants to lean toward choosing higher scores (due to social desirability bias) rather than

responding based on their true circumstances.

**Results** Table 30 presents the regression results predicting the self-reported personality in the human sample, the LLM-simulated personality, and the difference between human and LLM personalities, all based on the social desirability rating. Additionally, Figure 5 visually illustrates the trends through a regression plot.

It can be observed that when an item's social desirability rating is neutral (4), most LLMs tend to generate neutral responses (rated 3 on the scale). This result indicates that LLMs have learned to associate neutral social expectations with neutral opinions, consistent with previous research findings (Salecha et al., 2024). Regarding the mean difference between the human responses and the LLM responses, most differences show a positive correlation with the social desirability rating, and when the social desirability rating is neutral (4), the mean difference in most data is approximately 0.

In Figure 5, the line graph clearly illustrates that as the social desirability rating of the item increases, the average scores of both human responses and LLM responses rise accordingly. Moreover, as the social desirability rating increases, the average difference between human and LLM responses gradually widens, whereas this difference narrows when the social desirability scores are lower.

This trend indicates that when social desirability scores are high, the mean score of human responses tends to be higher than that of LLM responses. Conversely, when social desirability scores are low, the average score of human responses falls below that of LLM responses. This suggests that human responses are more strongly influenced by social desirability than those of LLMs.

Variable	Human	Persona GPT-4o	Shape GPT-4o	PSI GPT-4o	Persona Llama3	Shape Llama3	PSI Llama3
Intercept	1.15 (0.23)	1.90 (0.17)	2.33 (0.09)	2.13 (0.23)	0.71 (0.21)	2.25 (0.11)	2.08 (0.28)
Social desirability	0.46 (0.05)	0.23 (0.04)	0.12 (0.02)	0.23 (0.05)	0.47 (0.05)	0.14 (0.02)	0.20 (0.06)
$R^2$	.59	.39	.38	.27	.65	.37	.15
Predicted score at neutral point	2.99	2.82	2.81	3.05	2.59	2.81	2.88
Correlation	.77	.63	.62	.52	.80	.61	.39
Mean Score Difference							
Intercept	-	-0.76 (0.19)	-1.18 (0.19)	-0.99 (0.20)	0.44 (0.21)	-1.11 (0.23)	-0.93 (0.32)
Social desirability	-	0.23 (0.04)	0.34 (0.04)	0.23 (0.04)	-0.01 (0.05)	0.33 (0.05)	0.26 (0.07)
$R^2$	-	.34	.53	.32	.48	.42	.19
Predicted score at neutral point	-	0.16	0.18	-0.07	0.40	0.21	0.11
Correlation	-	.58	.73	.56	-.03	.65	.43

Table 30: Regression Analysis and Correlation of Social Desirability Ratings

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.

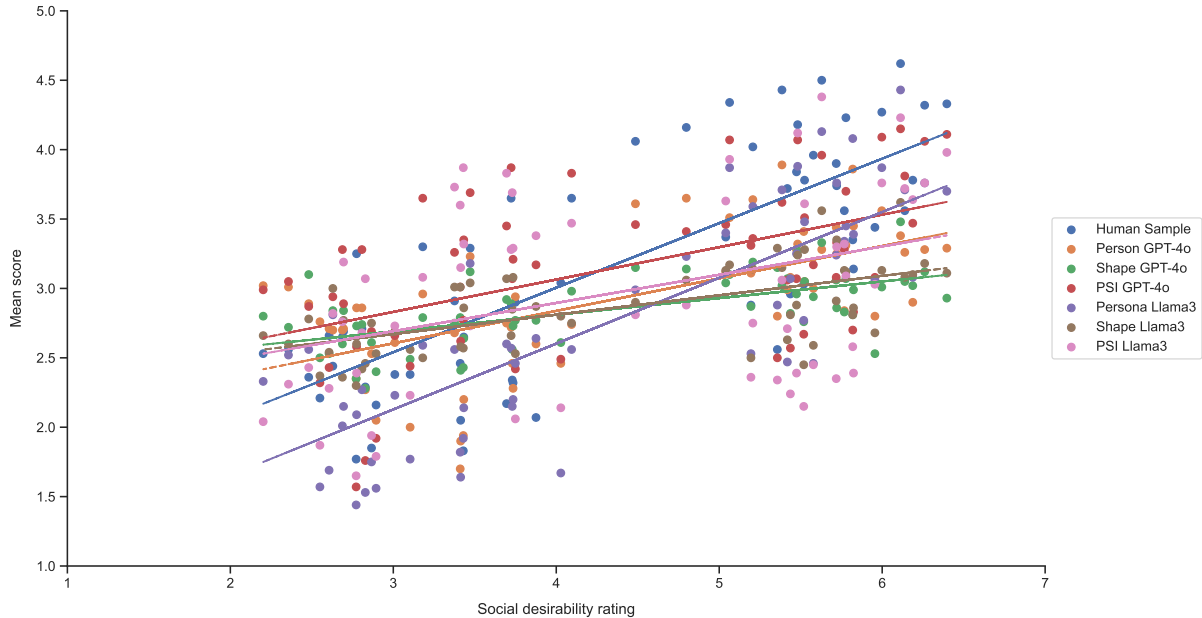


Figure 5: Regression Line Chart of Social Desirability Ratings

Note:  $n = 1,559$  for human responses,  $n = 297$  for Shape Llama3, and  $n = 300$  for other LLM responses. Some sample sizes are below 300, because certain generated data exceeded reasonable thresholds (1-5) for specific items and were excluded from the analysis.