

Rethinking Personality Assessment from Human-Agent Dialogues: Fewer Rounds May Be Better Than More

Anonymous ACL submission

Abstract

Personality assessment is essential for developing user-centered systems, playing a critical role across domains including hiring, education, and personalized system design. With the integration of conversational AI systems into daily life, automatically assessing human personality through natural language interaction has gradually gained more attention. However, existing personality assessment datasets based on natural language generally lack consideration of interactivity. Therefore, we propose Personality-1260, a Chinese dataset containing 1260 interaction rounds between humans and agents with different personalities, aiming to support research on personality assessment. Based on this dataset, we designed experiments to explore the effects of different interaction rounds and agent personalities on personality assessment. Results show that fewer interaction rounds perform better in most cases, and agents with different personalities stimulate different expressions of users' personalities. These findings provide guidance for the design of interactive personality assessment systems.

1 Introduction

Quantifying and benchmarking human behavior has always been an important topic in fields such as social science, philosophy, and psychology. As a core research direction, personality assessment not only helps reveal the internal mechanisms of individual behavioral patterns, thinking processes, and emotional responses, but also provides scientific evidence for mental health diagnosis (Widiger and Samuel, 2005), career planning (Tracey and Rounds, 1995), and educational method design (Bidjerano and Dai, 2007). With the emergence of chatbots and conversational AI systems becoming seamlessly integrated into daily life, automatically assessing human personality through natural language interaction has gradually gained more attention. From early dictionary-based tools

like LIWC (Pennebaker and King, 1999) to supervised learning model methods (Yang et al., 2021, 2023a), the rapid development of large language models (LLMs) provides unprecedented opportunities for dynamically capturing personality traits through natural language, such as PsyCoT (Yang et al., 2023b) and EERPD (Li et al., 2025).

Social Penetration Theory uses the "onion model" to describe personality (Altman and Taylor, 1973), which suggests that personality consists of multiple layers that are gradually revealed through interaction. However, current datasets for personality assessment through natural language lack consideration of interactivity. They mainly fall into two categories: one identifies personality traits from static texts like blogs (e.g. MBTI¹) and articles (e.g. Essays (Pennebaker and King, 1999)), which are easy to obtain but lack interactivity and struggle to reflect personality traits embedded in dynamic communication; the other uses manually annotated TV show or movie dialogues such as FriendsPersona (Jiang et al., 2020) and PersonalityEvd (Sun et al., 2024), providing interactive contexts but limited by acted and maybe exaggerated personalities, resulting in annotations lacking ecological validity in real environments. How to naturally and stably elicit comprehensive personality expressions at the language level in real interactive situations is key to effectively building datasets.

Media equation theory suggests that people unconsciously apply social rules when interacting with computers (Reeves and Nass, 1996). With advances in LLMs for human-agent interaction, combined with their excellent interactive capabilities in role-playing and personality simulation tasks (Shao et al., 2023; Chen et al., 2024; Jiang et al., 2024b), new opportunities have emerged. Compared to human-to-human dialogues, interactions with agents are more stable in long, multi-round

¹<https://www.kaggle.com/datasnaek/mbti-type>

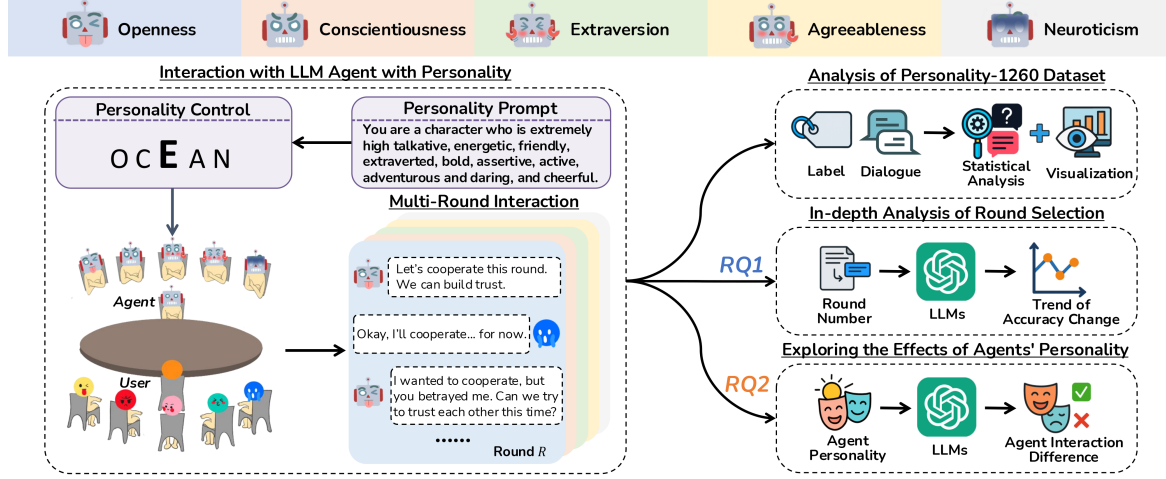


Figure 1: Workflow illustration of the paper. Left side shows the interaction scenario: controlling LLM agents through personality prompts to exhibit high Big Five dimension traits then interacting with users over multiple rounds. The right part is the core workflow, where we first analyzed Personality-1260 Dataset, then conducted experiments on the dimensions of interaction rounds and interacting agents to answer our research questions.

conversations (Guan et al., 2025), creating new chances for personality assessment. Leveraging these advantages, we developed five LLM agents using prompts based on the Big Five theory (Jiang et al., 2024a; Serapio-García et al., 2023) and designed game scenarios to constrain conversations and elicit personality expressions. Through these interactions with 42 real users, we constructed the **Personality-1260** dataset containing 1260 rounds of dialogues along with participants’ BFI-44 personality questionnaire results. This dataset helps study personality in human-agent interactions.

With Personality-1260 as data support, we explored personality assessment patterns in multi-round game scenarios between humans and agents with different personalities. In our research, we first validated the effectiveness of the dataset through statistical analysis and visualization. Then, based on these preliminary results, we compared the effectiveness of using different numbers of interaction rounds for personality assessment. Finally, we conducted further experiments by comparing interactions with agents having different personalities and their impact on assessment results. Building on these results, we aim to comprehensively evaluate personality assessment in human-agent interaction, focusing on the following research questions:

- **RQ1:** How much data do we need for effective personality assessment?
- **RQ2:** Does interacting with agents of different personalities influence personality assessment results?

2 Related Works

2.1 Personality

Personality refers to a stable structure formed by psychological and physiological systems within an individual, shaping and influencing their patterns of behavior, thoughts, and emotional responses (Allport, 1961). Psychologists have proposed various theories to understand personality, such as the Big Five (Briggs, 1992; Goldberg, 2013; De Raad, 2000), the Sixteen Personality Factors (16PF) (Cattell, 2001; Sells and Cattell, 1957), and the Myers-Briggs Type Indicator (MBTI) (Myers, 1962), all of which have seen extensive practical applications (Lounsbury et al., 2005). Among these theories, the Big Five is one of the most widely accepted (John et al., 2008), comprising Extraversion, Agreeableness, Conscientiousness, Neuroticism, and Openness. Each trait strongly correlates with specific behavioral tendencies (John, 1999). Beyond behavior, personality traits have also been found to correlate significantly with language use in communication (Hirsh and Peterson, 2009; Lee et al., 2007; Pennebaker and Graybeal, 2001; Pennebaker and King, 1999). Moreover, the Big Five have shown strong reliability and validity in cross-cultural studies (Gurven et al., 2013; Benet-Martínez and John, 1998). Therefore, this study adopts the Big Five framework as the foundation for analysis.

2.2 Automatic Personality Assessment

In recent years, automatic personality recognition has gained widespread attention due to its potential

to enhance personalized interactions (Qian et al., 2018; Zhang et al., 2018). Research in this field has evolved from analyzing language-based features to applying complex models. Early personality assessments primarily relied on linguistic features, such as the LIWC method, which predicted personality traits through language style and vocabulary usage (Francis and Booth, 1993). Later, traditional machine learning methods began to be applied in this field, such as the use of SVM (Cui and Qi, 2017) and XGBoost (Tadesse et al., 2018). However, these methods relied on manually extracted features, limiting their performance. The introduction of deep learning methods improved the accuracy of personality assessment. For example, Xue et al. combined hierarchical neural networks with the Inception variant to extract deep semantic features (Xue et al., 2018). The emergence of pre-trained models, such as BERT (Devlin et al., 2019), further enhanced performance. Keh et al. (Keh et al., 2019) and Jiang et al. (Jiang et al., 2020) used pre-trained models to extract features from posts and map user vectors to MBTI labels. TrigNet combined BERT initialization with a graph attention mechanism to integrate psycholinguistic knowledge (Yang et al., 2021). Despite these advances, these methods still face limitations in handling long texts.

Recently, LLMs have been applied to personality assessment. Some preliminary studies have used LLMs to decode personality traits from various forms of user-generated text (Peters et al., 2024; Peters and Matz, 2024; Zhang et al., 2024). Further research, such as that by Yang et al., combined Chain of Thought (CoT) with traditional personality questionnaires to predict personality traits (Yang et al., 2023b). Li et al. proposed a retrieval-augmented generation (RAG) framework, incorporating psychological knowledge of emotion regulation into LLM-based personality assessment (Li et al., 2025). Overall, while LLMs have shown promise in personality assessment, no study has yet explored the data requirements for LLM-based personality evaluation methods.

3 Dataset

3.1 Overview

Personality-1260 is a multi-round, multi-turn, dialogue-based dataset in Chinese (Fig. 2 shows the definitions of "round" and "turn") designed to assess personality by capturing authentic behaviors exhibited by human users during interactions

with agents of different personalities. The dataset includes Big Five personality dimension scale results from 42 participants (21 males, 21 females; $M = 22.07$, $SD = 2.32$) and records a total of 1,260 interaction rounds between humans and agents. Each round contained an average of 4.24 turns ($SD = 3.66$).

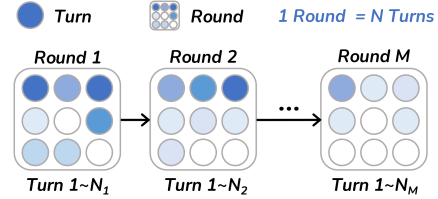


Figure 2: Definition of Round and Turn. The color of each turn represents interaction intensity (i.e., character count in dialogues). The number of interaction turns is not fixed but depends on users' discretion.

3.2 Data Collection Methods

3.2.1 Participants

We recruited 42 participants from a local university. All participants were fluent in the language used in the experiment. They all abstained from alcohol consumption, severe fatigue, drug use, or physical discomfort prior to the experiment. The study adhered to the principles of the Declaration of Helsinki and received approval from the Institutional Review Board. After being informed of general procedures and minimal risks, all participants provided written informed consent. To prevent bias like the social desirability effect, the specific purpose (i.e., personality trait assessment) was disclosed only after the experiment. During debriefing, participants were fully informed, received a US \$10 compensation, and were given the option to confirm or withdraw consent for data usage. Ultimately, all participants agreed to the use of their data for research purposes.

3.2.2 Experimental Environment Design

We developed a prototype system based on the Prisoner's Dilemma game as an interactive platform and deployed it on a personal computer (PC) (see Fig. 3). The Prisoner's Dilemma (Flood, 1958) is widely used in psychological experiments due to its effectiveness in simulating cooperative and defection behaviors in social contexts (Axelrod and Hamilton, 1981; Fehr and Gächter, 2002). Building on the traditional game mechanism, we introduced a natural dialogue exchange phase before the participants made their cooperation or defection de-

cisions. This addition aims to enhance interaction between the user and the agent, thereby simulating a more realistic interpersonal social scenario.

Prior research has shown that incorporating storylines can enhance immersion and engagement (Berson et al., 2018; Bouchard and Rizzo, 2019). Based on this, we designed a storyline to encourage participants to express their authentic selves during the game (see Appendix A.1). Notably, our storyline was not result-oriented (e.g., emphasizing score incentives or win-loss outcomes), but was designed to encourage users to fully express their true thoughts and behaviors. We deliberately minimized the emphasis on game mechanics to avoid interference with personality assessment (Jia et al., 2016) (for more details, please see Appendix A).

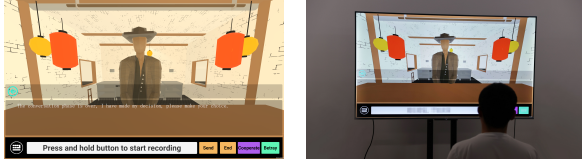


Figure 3: The system used in the experiment.

3.2.3 Experimental Procedure

Fig. 4 shows the the experimental procedure. It includes two phases: **Before Game**, **During Game**.

Before Game. Participants completed the Chinese version of the BFI-44 personality inventory (John and Srivastava, 1999), and familiarized themselves with the system operation. They were then instructed to carefully read the storyline described in Section 3.2.2, along with the rules of the Prisoner’s Dilemma game, where players can choose to cooperate or defect—cooperation benefits both sides, but defection may yield greater advantage for one player (see Appendix A.2 for details).

During Game. Participants interacted with five LLM agents that exhibited the most significant characteristics (highest scores) on each dimension of the Big Five: Openness (O), Conscientiousness (C), Extraversion (E), Agreeableness (A), and Neuroticism (N). **The interaction sequence was randomized across participants.** Interaction with each agent consisted of six rounds, each comprising a dialogue phase and a decision phase:

- **Dialogue Phase:** Participants could communicate freely with the agent via voice or text to influence its decisions.

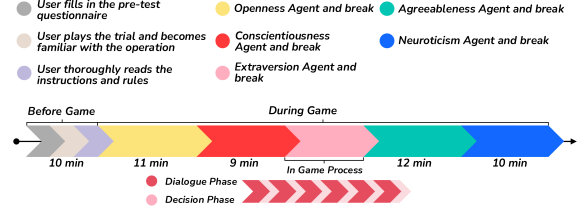
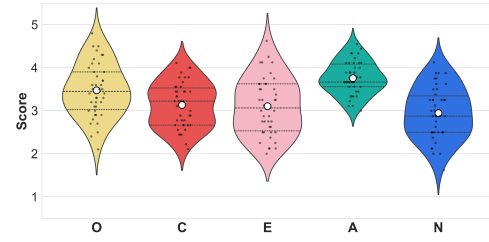
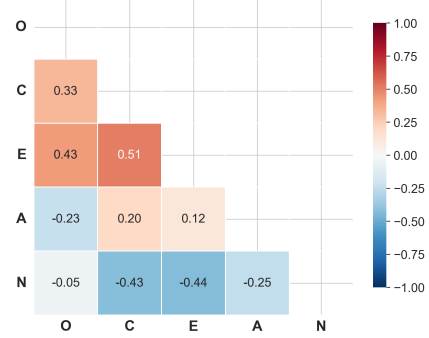


Figure 4: Overview of the experimental procedure. In experiment, the sequence of agent is randomized.



(a) Distribution of personality trait scores across all participants. White dots indicate mean values.



(b) Correlation matrix between Big Five personality dimensions.

Figure 5: Label distribution of Personality-1260 dataset: (a) violin plot illustrating distribution patterns, (b) correlation heatmap revealing relationships between traits.

- **Decision Phase:** Both parties independently chose “cooperate” or “defect.”

The number of game rounds was determined based on small-scale user testing during development, ensuring interactions lasted approximately 10 minutes to maintain engagement without causing fatigue. The number of dialogue exchanges (turns) per round was at the participant’s, and participants could end the dialogue at any time.

3.3 Dataset Statistics

3.3.1 Label Statistics

We visualized the label distributions of the Personality-1260 dataset. As shown in Fig. 5a, Agreeableness had the highest average score ($M = 3.76$, $SD = 0.38$), followed by Openness ($M = 3.47$, $SD = 0.62$). In contrast, Conscientiousness ($M = 3.13$, $SD = 0.53$), Extraversion ($M =$

3.10, $SD = 0.67$), and Neuroticism ($M = 2.95$, $SD = 0.61$) had progressively lower scores. This ranking aligns with the findings of Zhang et al. (Zhang et al., 2022), supporting the validity of our dataset. Additionally, Agreeableness scores were most concentrated (3.0–4.56). In comparison, Extraversion showed the greatest variability ($SD = 0.67$), while Neuroticism had the widest score range (1.63–4.13). These results suggest substantial individual differences in these two traits, reflecting the diversity of the dataset.

The correlation heatmap in Fig. 5b highlights five significant correlations ($|r| \geq 0.3$). A relatively strong positive correlation was observed between Extraversion and Conscientiousness ($r = 0.51$). Although this correlation was higher than in previous studies (Zhao and Seibert, 2006), it aligns with findings indicating that Extraversion and Conscientiousness often jointly predict positive life outcomes (Soto and John, 2017; Vella, 2024). Additionally, moderate positive correlations were found between Extraversion and Openness ($r = 0.43$), and between Openness and Conscientiousness ($r = 0.33$), consistent with Liu et al. (Liu and Campbell, 2017). Meanwhile, significant negative correlations appeared between Neuroticism and Extraversion ($r = -0.44$), as well as between Neuroticism and Conscientiousness ($r = -0.43$). These negative correlations align with previous Big Five personality research (Van der Linden et al., 2010), further confirming the validity of our dataset.

3.3.2 Dialogue Statistics

The Personality-1260 dataset includes multiple rounds of interactions between users and an agent. Therefore, we further analyzed how user-agent interactions change over time. Specifically, we visualized the average number of turns per round and the average number of characters generated by users per round. As shown in Fig. 6, clear trends emerged during the six rounds of interaction. The average number of turns per round was highest in the first round (approximately 6.0 turns) but showed a clear decrease in the second round to around 4.0 turns, then remained relatively stable between 3.7 and 4.1 turns in subsequent rounds. A similar declining trend was observed for the average number of characters generated per round by users. This gradual reduction in linguistic output may indicate a decrease in user engagement as the interactions progressed.

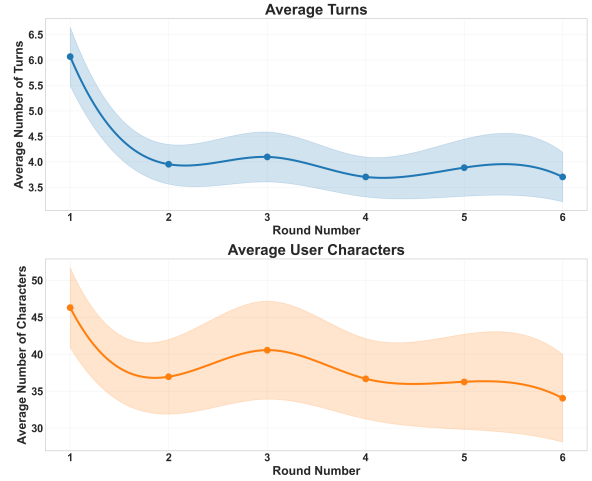


Figure 6: Temporal changes across six rounds: (a) average number of turns per round; (b) average number of characters generated by the user per round. Shaded areas represent 95% confidence intervals.

4 Experiments

Section 4.1 presents our experimental task formulation and implementation details. In Section 4.2, we outline the evaluation metrics, followed by our experimental design in Section 4.3. The corresponding experimental results are detailed across Sections 4.4, 4.5, and 4.6. Drawing from these findings, Section 4.7 offers three design recommendations for interactive personality assessment systems.

4.1 Experimental Setup

4.1.1 Task Formulation

There are five Big-Five personality dimensions $BF = [bf_1, bf_2, \dots, bf_5]$. Each dialogue D consists of interactions between a user U and an agent A . The dialogue D consists of 6 rounds $R = [r_1, r_2, \dots, r_6]$, where each round $r_i = [t_{i,1}, t_{i,2}, \dots, t_{i,n_i}]$ consists of multiple turns of conversation between the user and agent. This task aims to predict a score vector $P = [p_1, p_2, \dots, p_5]$ by minimizing the distributional difference between P and the ground-truth personality vector BF , and to provide supporting evidence $E = [e_1, e_2, \dots, e_5]$, where each e_j contains specific dialogue excerpts justifying the assigned score p_j .

4.1.2 Implementation Details

We implemented our experiment pipeline in Python using the OpenAI/Deepseek API. All experiments were conducted on a MacBook Pro with an M4 Pro chip. We set the temperature to 0 to get a reliable rather than innovative output. All experiments were run 3 times and the average values were taken.

	GPT-4.1-Nano						GPT-4.1					
Rounds	O	C	E	A	N	AVG	O	C	E	A	N	AVG
1	0.622	0.556	0.649	0.457	0.819	0.621	0.652	0.583	0.631	0.676	0.721	0.653
1-2	0.610	0.615	0.672	0.493	0.833	0.644	0.649	0.609	0.617	0.628	0.737	0.648
1-3	0.607	0.593	0.679	0.500	0.851	0.646	0.643	0.631	0.609	0.640	0.725	0.650
1-4	0.601	0.624	0.652	0.507	0.845	0.646	0.640	0.639	0.602	0.629	0.718	0.646
1-5	0.649	0.655	0.675	0.519	0.825	0.664	0.637	0.653	0.612	0.625	0.717	0.649
1-6	0.579	0.662	0.694	0.515	0.860	0.662	0.652	0.676	0.612	0.613	0.717	0.654
	GPT-4.1-Mini						DeepSeek-V3					
1	0.633	0.612	0.657	0.589	1.076	0.713	0.960	0.681	0.842	0.890	1.196	0.914
1-2	0.658	0.649	0.619	0.552	1.036	0.703	0.970	0.681	0.860	0.887	1.204	0.920
1-3	0.679	0.671	0.629	0.580	1.002	0.712	0.965	0.714	0.819	0.941	1.202	0.928
1-4	0.677	0.703	0.635	0.567	0.955	0.707	0.952	0.726	0.831	0.989	1.115	0.923
1-5	0.700	0.712	0.634	0.600	0.973	0.724	0.936	0.713	0.790	0.954	1.110	0.901
1-6	0.720	0.720	0.649	0.606	0.959	0.731	0.941	0.754	0.812	1.020	1.136	0.933

Table 1: MAE scores of different models across cumulative interaction rounds. Bolded values indicate the best performance among different cumulative round combinations. Columns O, C, E, A, N represent the MAE for the five dimensions of the Big Five model, while the AVG column represents the average value across all five dimensions.

4.2 Evaluation Metrics

To quantitatively assess the accuracy of our personality assessment results, we use the Mean Absolute Error (MAE) as the evaluation metric. For personality assessment on a standardized scale, MAE provides an intuitive measure of prediction accuracy. The MAE is calculated as:

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (1)$$

4.3 Experiment Design

Our experimental design includes two main dimensions (as shown in Fig. 7): Interaction Round Dimension and Interaction Agent Dimension. To answer our two research questions, "RQ1: How much data do we need for effective personality assessment?" and "RQ2: Does interacting with agents of different personalities influence personality assessment results?", we designed experiments on these two dimensions.

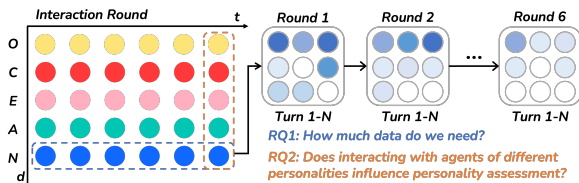


Figure 7: Illustration of experiment design.

Interaction Round Dimension: Multiple interaction rounds, each round includes 1 to N turns.

Interaction Agent Dimension: Different agents exhibiting high levels of traits in the Big Five.

4.4 In-depth Analysis of Round Selection

To answer the first research question regarding data requirements for personality assessment, we evaluated four state-of-the-art large language models: GPT-4.1-Nano, GPT-4.1-Mini, GPT-4.1 and DeepSeek-V3. Table 1 presents the Mean Absolute Error (MAE) scores for each of the Big Five personality dimensions (Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism) across different interaction rounds, where lower MAE values indicate better assessment accuracy.

Finding 1: The optimal data requirement for personality assessment appears to be 1-2 rounds of interaction. We conducted paired t-tests between all rounds (for example, comparing data from Round 1 with Round 1-6) and extracted round pairs with significant differences. Results show that in most cases, using data from the first two rounds of interaction for assessment produces the lowest error rates (see in Tables 10, 19, 28, 37). Contrary to intuitive expectations, in most cases, extending the number of rounds yields decreases in performance or no improvement. This finding has important practical implications for personality assessment

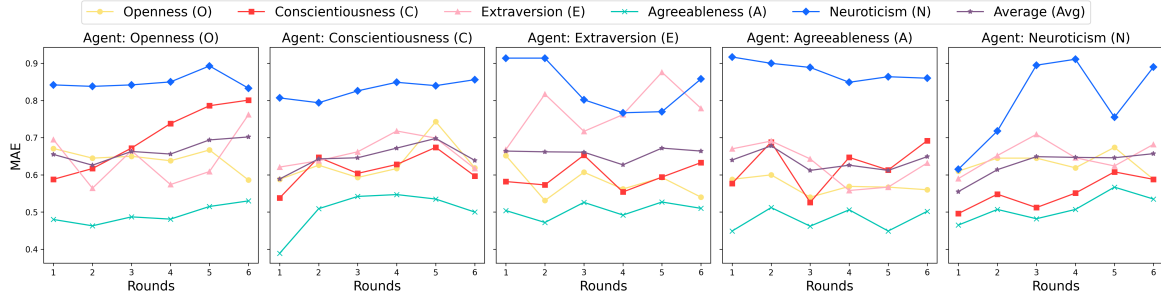


Figure 8: MAE scores across six interaction rounds between human participants and LLM agents. Each panel represents interactions with an agent exhibiting high levels of different personality dimensions. Colored lines represent MAE values for different personality dimensions.

system design, indicating that brief, targeted interactions may be more effective for personality assessment tasks than lengthy conversations.

Finding 2: Different personality dimensions show varying sensitivity to interaction duration.

We observed that Conscientiousness exhibits a significant or near-significant increasing trend in error across all models (see Table 3, 12, 21, 30). This may be because the characteristics associated with Conscientiousness (such as organization, discipline, and attention to detail) tend to become diluted or even contradictory as the conversation expands to cover more topics. In contrast, Extraversion and Openness, except in the GPT-4.1-Mini model, do not show significant trends of increasing or decreasing error. According to Table 1, the best measurement results for Openness mostly appear at the end of the interactions, which may indicate that as the interaction progresses, the assessment of a user’s Openness becomes more accurate.

Finding 3: Neuroticism is difficult to assess accurately through dialogue analysis.

Across all models and interaction lengths, Neuroticism consistently shows the highest MAE scores among the five dimensions. DeepSeek-V3’s error rates for Neuroticism are highest, and even for GPT-4.1, which performs best in this dimension, Neuroticism error rates exceed those of other dimensions. This may be because emotional stability traits are inherently more difficult to detect from text-based interactions, while other dimensions manifest as more explicit behavioral descriptions. GPT-4.1 may achieve relatively better results due to its excellent performance in emotion perception.

Finding 4: Models with larger parameter counts may demonstrate better stability in assessments.

GPT-4.1 and DeepSeek-V3 show greater stability

in assessing Extraversion, Agreeableness, and Neuroticism compared to GPT-4.1-Nano and GPT-4.1-Mini, even though they may sometimes have larger errors than smaller parameter models.

4.5 Exploring the Effects of Agents’ Personality

To address our second research question (RQ2: "Does interacting with agents of different personalities influence personality assessment results?"), we conducted experiments using GPT-4.1-Nano, which performed best in our task. Fig. 8 presents the MAE scores across personality dimensions when interacting with agents exhibiting high levels of different personality dimensions.

Finding 1: Agent personality influences the accuracy of personality dimension assessment.

Most notably, in the condition of interacting with a neuroticism agent, the error in the Neuroticism dimension in the first round is significantly lower than when interacting with agents of other personalities. Mann-Whitney U tests revealed significant differences between Neuroticism agents and Agreeableness ($U = 614.5, p = 0.008$), Extraversion ($U = 610.0, p = 0.0075$), and Openness agents ($U = 627.0, p = 0.011$), with a marginally significant difference compared to Conscientiousness agents ($U = 722.5, p = 0.077$). This may be because Neuroticism agent produces stronger stimuli for users in the first round of interaction, evoking manifestations of their Neuroticism traits, while users show adaptability in subsequent rounds.

Similarly, in the first round of interaction with high Conscientiousness agents, optimal assessment of user Agreeableness was achieved. We computed Cohen’s d for Conscientiousness versus each other agent type, with all effect sizes falling in the small ($|d| \approx 0.2$) to small-to-medium ($|d| \approx 0.3$) range

(C vs. A: $d = -0.21$; C vs. E: $d = -0.30$; C vs. O: $d = -0.31$; and C vs. N: $d = -0.20$). This may be due to the organizational, disciplined, and polite characteristics of Conscientiousness agents also evoking manifestations of Agreeableness traits.

Finding 2: Specific trait agents can be deployed when assessing specific dimensions. As noted in Finding 1, interactions with agents of different traits have varying effects on assessing specific dimensions. When assessment systems need to focus on specific personality dimensions, the corresponding agent type should be carefully selected. For example, when assessing Agreeableness, data from the first round of interaction with a high Conscientiousness agent may be chosen; when assessing Neuroticism, data from the first round of interaction with a high Neuroticism agent should be used.

Finding 3: Assessment of the Openness dimension can benefit from appropriate attention to interaction duration. We found that when assessing Openness, interactions with high Agreeableness, high Extraversion, and high Openness agents show decreasing errors as interaction duration increases, which is consistent with Table 1.

4.6 Comparison with Human Annotators

To better validate our findings, we recruited four senior PhD students in psychology to annotate the content in our dataset. We used Intraclass Correlation Coefficient (ICC) analysis and Friedman tests to evaluate the rating consistency and differences among the four annotators. Results showed that despite high overall consistency ($ICC \geq 0.60$), significant systematic differences still existed among annotator ratings across the five dimensions ($p \leq 0.001$), indicating annotators generally agreed on which users had stronger or weaker traits but differed in their overall rating tendencies (for more details, please see Appendix G.1.4 and G).

Additionally, we observed a trend in Table 40 that aligns with Table 1: In most cases, extending the number of rounds yields decreases in performance or no improvement (Tables 41, 42, 43, and 44 show linear tests of error trends and round pairs with significant differences ($p < 0.05$)). Furthermore, we found that the evaluation results from LLMs were comparable to those from human evaluators. Overall, the human annotation results support our experimental findings and highlight the importance of including real user labels in the dataset.

4.7 Design Recommendations

The above two experiments reveal several important findings, such as "more" does not equal "better." Experiment One indicates that increasing interaction rounds may actually reduce assessment accuracy, with the optimal data volume typically being 1-2 rounds of interaction. Experiment Two demonstrates that the importance of specific agent-dimension matching may exceed the data volume.

These findings provide several recommendations for interactive personality assessment systems:

- Optimizing specific interaction quality (e.g., appropriate agent–dimension matching) is more important than simply increasing the number of interaction rounds.
- Different approaches may be needed for assessing different personality dimensions. For example, when evaluating Openness, we should consider the fragmented features which users exhibit in long-term interactions.
- The complex effects of the interaction environment and the number of interaction rounds should be considered when designing personality assessment systems. For example, when assessing Neuroticism, we could use first-round interaction data with the agent exhibiting strong Neuroticism traits.

5 Conclusion

In this study, we focused on personality assessment in human-agent interaction and introduced Personality-1260, addressing the gap in existing datasets that lack either interactivity or authentic user labels. We validated this dataset’s effectiveness through statistical analysis and visualization. Based on Personality-1260, we experimentally explored how different interaction rounds and agent personalities influence personality assessment. Contrary to intuition, our results demonstrated that in most cases, extending the number of rounds either decreases performance or yields no improvement. Additionally, we found that the interacting agent’s personality influences the accuracy of personality assessment. Based on these experimental findings, we proposed three design recommendations for interactive personality assessment systems. We hope these insights can provide guidance for the future design of interactive personality assessment systems.

Limitations

There are several limitations of our Personality-1260 dataset and experiments.

First, our dataset is in Chinese. Although the Big Five personality traits have been validated to have good generalizability across cultural samples, the ideal scenario would still be to build multilingual datasets to support personality assessment across different cultures.

Second, our participant demographics are not sufficiently diverse, as all participants came from one university. However, by analyzing these participants' Big Five questionnaire results, we found a high degree of overlap with distributions from previous studies with broader participant demographics, which also validates the effectiveness of our dataset.

Third, compared to the two existing types of datasets (those based on static texts like writing/social media, and those manually annotated from TV shows/movies), our dataset is not large. However, we have filled the gap between them - Personality-1260 has both dynamic interactivity and real personality labels from users. Moreover, it is sufficient in diversity and depth to support meaningful analysis. We plan to further expand the dataset in the future.

Finally, this study mainly focuses on closed-source GPT series models and a small number of open-source models. We had experimented with the open-source Qwen-2.5-plus, where the average MAE score for each dimension was around 2, indicating that the assessment error was extremely large, lacked reference value, and was not suitable for experimental analysis. Because the performance of Qwen-2.5-plus was not good and given budget constraints, we conducted experiments on GPT-4.1-nano, GPT-4.1-mini, GPT-4.1, and deepseek-v3.

Ethics Statements

This study strictly adheres to the ACL Code of Ethics for human experiments and has received approval from the Institutional Review Board (IRB). The experiment lasted approximately one hour, with each participant receiving a compensation of \$10, which constitutes a fair and reasonable hourly wage in the local area. To avoid biases such as the social desirability effect, the specific purpose of the study (personality trait assessment) was only disclosed after the experiment. During the debriefing

session, participants were fully informed and given the option to confirm or withdraw their consent for data usage. Ultimately, all participants agreed to the use of their data for research purposes and provided written informed consent.

With the increasing prevalence of AI dialogue systems in daily life, massive amounts of data have become available for interactive personality assessment. However, this technological advancement also comes with potential risks, and we must remain vigilant against its possible use for harmful purposes targeting individuals, groups, or society. These risks include unauthorized personality analysis, targeted manipulation, and privacy violations, which are particularly severe when users are unaware.

Based on Responsible AI principles, we have implemented multiple protective measures. Regarding privacy protection, we strictly adhere to data confidentiality principles, ensuring that all personal data is secure and used solely for research purposes. In terms of transparency, we have disclosed the experimental prompts in the paper's appendix, enhancing the reproducibility of our research. During the personality assessment process, we required LLMs to provide evidence-based, traceable results, ensuring the reliability and fairness of the assessments. We strongly advocate the research community to maintain high vigilance regarding data and privacy security, ensuring that users are fully informed and participate voluntarily, while clearly defining the purposes of data collection and strictly limiting its scope.

Our research aims to analyze the key factors affecting interactive personality assessment, to support the design of better personality assessment systems that help users gain deeper self-understanding and subsequently support their career planning and personal development. Through rigorous ethical review and informed consent procedures, we strive to balance technological innovation with ethical responsibility, ensuring that advances in AI-assisted personality assessment truly benefit individuals and society without compromising personal rights or well-being.

References

- Gordon W Allport. 1961. Pattern and growth in personality.
- Irwin Altman and Dalmis A Taylor. 1973. *Social pene-*

690	<i>tration: The development of interpersonal relationships.</i> Holt, Rinehart & Winston.	ME Francis and Roger J Booth. 1993. Linguistic inquiry and word count. <i>Southern Methodist University: Dallas, TX, USA.</i>	742
691			743
692	Robert Axelrod and William D Hamilton. 1981. The evolution of cooperation. <i>science</i> , 211(4489):1390–1396.	Lewis R Goldberg. 2013. An alternative “description of personality”: The big-five factor structure. In <i>Personality and Personality Disorders</i> , pages 34–47. Routledge.	745
693			746
694			747
695	Verónica Benet-Martínez and Oliver P John. 1998. Los cinco grandes across cultures and ethnic groups: Multitrait-multimethod analyses of the big five in spanish and english. <i>Journal of personality and social psychology</i> , 75(3):729.	Shengyue Guan, Haoyi Xiong, Jindong Wang, Jiang Bian, Bin Zhu, and Jian-guang Lou. 2025. Evaluating llm-based agents for multi-turn conversations: A survey. <i>arXiv preprint arXiv:2503.22458</i> .	748
696			749
697			750
698			751
699			752
700	Ilene R Berson, Michael J Berson, Amy M Carnes, and Claudia R Wiedeman. 2018. Excursion into empathy: exploring prejudice with virtual reality. <i>Social Education</i> , 82(2):96–100.	Michael Gurven, Christopher Von Rueden, Maxim Massenkov, Hillard Kaplan, and Marino Lero Vie. 2013. How universal is the big five? testing the five-factor model of personality variation among forager-farmers in the bolivian amazon. <i>Journal of personality and social psychology</i> , 104(2):354.	753
701			754
702			755
703			756
704	Temi Bidjerano and David Yun Dai. 2007. The relationship between the big-five model of personality and self-regulated learning strategies. <i>Learning and individual differences</i> , 17(1):69–81.	Jacob B. Hirsh and Jordan B. Peterson. 2009. <i>Personality and language use in self-narratives</i> . <i>Journal of Research in Personality</i> , page 524–527.	757
705			758
706			
707			759
708	Stéphane Bouchard and A Rizzo. 2019. <i>Virtual reality for psychological and neurocognitive interventions</i> . Springer.	Yuan Jia, Bin Xu, Yamini Karanam, and Stephen Voids. 2016. Personality-targeted gamification: a survey study on personality traits and motivational affordances. In <i>Proceedings of the 2016 CHI conference on human factors in computing systems</i> , pages 2001–2013.	760
709			761
710			762
711	Stephen R. Briggs. 1992. <i>Assessing the five-factor model of personality description</i> . <i>Journal of Personality</i> , 60:253–293.	Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2024a. Evaluating and inducing personality in pre-trained language models. <i>Advances in Neural Information Processing Systems</i> , 36.	763
712			764
713			765
714	Heather EP Cattell. 2001. The sixteen personality factor (16pf) questionnaire. In <i>Understanding psychological assessment</i> , pages 187–215. Springer.	Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2024b. <i>PersonaLLM: Investigating the ability of large language models to express personality traits</i> . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 3605–3627, Mexico City, Mexico. Association for Computational Linguistics.	766
715			767
716			768
717	Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024. <i>From persona to personalization: A survey on role-playing language agents</i> . <i>Transactions on Machine Learning Research</i> . Survey Certification.	Hang Jiang, Xianzhe Zhang, and Jinho D Choi. 2020. Automatic text-based personality recognition on monologues and multiparty dialogues using attentive networks and contextual embeddings (student abstract). In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 13821–13822.	769
718			770
719			771
720			772
721			773
722			774
723			775
724	Brandon Cui and Calvin Qi. 2017. Survey analysis of machine learning methods for natural language processing for mbti personality type prediction. <i>Final Report Stanford University</i> .	O John. 1999. The big five trait taxonomy: History, measurement, and theoretical perspectives. <i>Handbook of personality/Guilford</i> .	776
725			777
726			778
727			779
728	Boele De Raad. 2000. <i>The big five personality factors: the psycholexical approach to personality</i> . Hogrefe & Huber Publishers.	Oliver P John, Laura P Naumann, and Christopher J Soto. 2008. Paradigm shift to the integrative big five trait taxonomy. <i>Handbook of personality: Theory and research</i> , 3(2):114–158.	780
729			781
730			782
731	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)</i> , pages 4171–4186.	Oliver P. John and Sanjay Srivastava. 1999. <i>Handbook of Personality: Theory and Research</i> , 2nd edition. Guilford Press, New York. Chinese edition:	783
732			784
733			785
734			786
735			787
736			788
737			789
738	Ernst Fehr and Simon Gächter. 2002. Altruistic punishment in humans. <i>Nature</i> , 415(6868):137–140.		790
739			791
740	Merrill M Flood. 1958. Some experimental games. <i>Management Science</i> , 5(1):5–26.		792
741			793
			794
			795
			796

797	Lawrence A. Pervin, Oliver P. John, 2003:135–184.	851
798	(Chinese BFI-44 printed on p.176 of the Chinese	852
799	edition).	853
800	Sedrick Scott Keh, I Cheng, and 1 others. 2019. Myers-	854
801	briggs personality classification and personality-	855
802	specific language generation using pre-trained lan-	856
803	guage models. <i>arXiv preprint arXiv:1907.06333</i> .	857
804	Chang H. Lee, Kyungil Kim, Young Seok Seo, and	858
805	Cindy K. Chung. 2007. The relations between per-	
806	sonality and language use . <i>The Journal of General</i>	
807	<i>Psychology</i> , 134:405–413.	
808	Zheng Li, Dawei Zhu, Qilong Ma, Weimin Xiong, and	
809	Sujian Li. 2025. EERPD: Leveraging emotion and	
810	emotion regulation for improving personality detec-	
811	tion . In <i>Proceedings of the 31st International Con-</i>	
812	<i>ference on Computational Linguistics</i> , pages 7721–	
813	7734, Abu Dhabi, UAE. Association for Computa-	
814	tional Linguistics.	
815	Dong Liu and W Keith Campbell. 2017. The big five	
816	personality traits, big two metatraits and social media:	
817	A meta-analysis. <i>Journal of Research in Personality</i> ,	
818	70:229–240.	
819	John W Lounsbury, Teresa Hutchens, and James M	
820	Loveland. 2005. An investigation of big five person-	
821	ality traits and career decidedness among early and	
822	middle adolescents. <i>Journal of career assessment</i> ,	
823	13(1):25–39.	
824	IB Myers. 1962. The myers-briggs type indicator. <i>Edu-</i>	
825	<i>cational Testing Service/Princeton</i> .	
826	James W Pennebaker and Anna Graybeal. 2001. Pat-	
827	terns of natural language use: Disclosure, personality,	
828	and social integration. <i>Current Directions in Psycho-</i>	
829	<i>logical Science</i> , 10(3):90–93.	
830	James W Pennebaker and Laura A King. 1999. Lin-	
831	guistic styles: language use as an individual differ-	
832	ence. <i>Journal of personality and social psychology</i> ,	
833	77(6):1296.	
834	Heinrich Peters, Moran Cerf, and Sandra C Matz.	
835	2024. Large language models can infer personal-	
836	ity from free-form user interactions. <i>arXiv preprint</i>	
837	<i>arXiv:2405.13052</i> .	
838	Heinrich Peters and Sandra Matz. 2024. Large language	
839	models can infer psychological dispositions of social	
840	media users. <i>PNAS Nexus</i> , 3(6):pgae231.	
841	Qiao Qian, Minlie Huang, Haizhou Zhao, Jingfang	
842	Xu, and Xiaoyan Zhu. 2018. Assigning personal-	
843	ity/profile to a chatting machine for coherent con-	
844	versation generation . In <i>Proceedings of the Twenty-</i>	
845	<i>Seventh International Joint Conference on Artificial</i>	
846	<i>Intelligence</i> .	
847	Byron Reeves and Clifford Nass. 1996. The media	
848	equation: How people treat computers, television,	
849	and new media like real people. <i>Cambridge, UK</i> ,	
850	10(10):19–36.	
	S. B. Sells and Raymond B. Cattell. 1957. Personal-	851
	ity and motivation structure and measurement . <i>The</i>	852
	<i>American Journal of Psychology</i> , page 620.	853
	Greg Serapio-García, Mustafa Safdari, Clément Crepy,	854
	Luning Sun, Stephen Fitz, Peter Romero, Marwa	855
	Abdulhai, Aleksandra Faust, and Maja Matarić. 2023.	856
	Personality traits in large language models. <i>arXiv</i>	857
	<i>preprint arXiv:2307.00184</i> .	858
	Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu.	859
	2023. Character-llm: A trainable agent for role-	860
	playing. <i>arXiv preprint arXiv:2310.10158</i> .	861
	Christopher J Soto and Oliver P John. 2017. The next	862
	big five inventory (bfi-2): Developing and assess-	863
	ing a hierarchical model with 15 facets to enhance	864
	bandwidth, fidelity, and predictive power. <i>Journal of</i>	865
	<i>personality and social psychology</i> , 113(1):117.	866
	Lei Sun, Jinming Zhao, and Qin Jin. 2024. Reveal-	867
	ing personality traits: A new benchmark dataset	868
	for explainable personality recognition on dialogues .	869
	In <i>Proceedings of the 2024 Conference on Empiri-</i>	870
	<i>cal Methods in Natural Language Processing</i> , pages	871
	19988–20002, Miami, Florida, USA. Association for	872
	Computational Linguistics.	873
	Michael M Tadesse, Hongfei Lin, Bo Xu, and Liang	874
	Yang. 2018. Personality predictions based on user be-	875
	havior on the facebook social media platform. <i>IEEE</i>	876
	<i>Access</i> , 6:61959–61969.	877
	Terence JG Tracey and James Rounds. 1995. The arbi-	878
	trary nature of holland’s riasec types: A concentric-	879
	circles structure. <i>Journal of Counseling Psychology</i> ,	880
	42(4):431.	881
	Dimitri Van der Linden, Jan Te Nijenhuis, and Arnold B	882
	Bakker. 2010. The general factor of personality:	883
	A meta-analysis of big five intercorrelations and a	884
	criterion-related validity study. <i>Journal of research</i>	885
	<i>in personality</i> , 44(3):315–327.	886
	Melchior Vella. 2024. The relationship between the	887
	big five personality traits and earnings: Evidence	888
	from a meta-analysis. <i>Bulletin of Economic Research</i> ,	889
	76(3):685–712.	890
	Thomas A Widiger and Douglas B Samuel. 2005.	891
	Evidence-based assessment of personality disorders.	892
	<i>Psychological Assessment</i> , 17(3):278.	893
	Di Xue, Lifa Wu, Zheng Hong, Shize Guo, Liang Gao,	894
	Zhiyong Wu, Xiaofeng Zhong, and Jianshan Sun.	895
	2018. Deep learning-based personality recognition	896
	from text posts of online social networks. <i>Applied</i>	897
	<i>Intelligence</i> , 48(11):4232–4246.	898
	Tao Yang, Jinghao Deng, Xiaojun Quan, and Qifan	899
	Wang. 2023a. Orders are unwanted: dynamic deep	900
	graph convolutional network for personality detec-	901
	tion. In <i>Proceedings of the AAAI Conference on Arti-</i>	902
	<i>ficial Intelligence</i> , volume 37, pages 13896–13904.	903

Tao Yang, Tianyuan Shi, Fanqi Wan, Xiaojun Quan, Qifan Wang, Bingzhe Wu, and Jiaxiang Wu. 2023b. Psycot: Psychological questionnaire as powerful chain-of-thought for personality detection. *arXiv preprint arXiv:2310.20256*.

Tao Yang, Feifan Yang, Haolan Ouyang, and Xiaojun Quan. 2021. [Psycholinguistic tripartite graph network for personality detection](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4229–4239, Online. Association for Computational Linguistics.

Bo Zhang, Yi Ming Li, Jian Li, Jing Luo, Yonghao Ye, Lu Yin, Zhuosheng Chen, Christopher J Soto, and Oliver P John. 2022. The big five inventory–2 in china: A comprehensive psychometric evaluation in four diverse samples. *Assessment*, 29(6):1262–1284.

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?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Tianyi Zhang, Antonis Koutsoumpis, Janneke K Oostrom, Djurre Holtrop, Sina Ghassemi, and Reinout E de Vries. 2024. Can large language models assess personality from asynchronous video interviews? a comprehensive evaluation of validity, reliability, fairness, and rating patterns. *IEEE Transactions on Affective Computing*.

Hao Zhao and Scott E Seibert. 2006. The big five personality dimensions and entrepreneurial status: a meta-analytical review. *Journal of applied psychology*, 91(2):259.

A Appendix

A.1 Storyline

In a uniquely styled Eastern restaurant, you find yourself standing at the bar, facing a mysterious cowboy. He’s wearing a wide-brimmed hat and an old-fashioned trench coat, seemingly waiting for your next move. This isn’t just a casual encounter; it’s a crucial game. The room is simply decorated but carries an air of deep mystery. Red lanterns sway gently on either side, casting a warm orange glow on your face.

You’ve been selected by a secret organization to participate in this highly challenging game. The organization has informed you that the outcome of this game will have profound implications for its future, but they haven’t told you what result would be favorable. They only emphasized one thing—you must act according to your true thoughts and show your most authentic self. Your opponents aren’t just one person; they may look the same, but each one is different.

Remember, this is not just a game, but also an opportunity for self-discovery and expression. Regardless of the final outcome, as long as you stay true to your heart, there will be no regrets. Now, the game is about to begin—are you ready to face the challenge?

A.2 Game Rules

To help you better engage in this game, here are the rules:

1. Each round consists of two phases: the Dialogue Phase and the Decision Phase.
2. During the Dialogue Phase, you and your opponent can freely converse to influence each other’s decisions, such as building trust or making threats.
3. In the Decision Phase, both you and your opponent must independently choose either "Cooperate" or "Defect," which is the only way to interact with the game system.

4. If both players choose to cooperate, you will each earn 2 points.
5. If one player chooses to cooperate while the other chooses to defect, the defector will earn 3 points, and the co-operator will receive 0 points.
6. If both players choose to defect, you will each receive 0 points.

Are you ready to enter this unknown territory and face the challenge?

A.3 Personality Control

Extraversion: You are a character who is extremely high in talkativeness, energy, friendliness, extraversion, boldness, assertiveness, activeness, adventurousness, daringness, and cheerfulness.

Agreeableness: You are a character who is extremely high in altruism, cooperativeness, trust, morality, honesty, kindness, generosity, humbleness, sympathy, unselfishness, and agreeableness.

Conscientiousness: You are a character who is extremely high in responsibility, hardworkingness, self-efficacy, orderliness, self-discipline, practicality, thriftiness, organization, conscientiousness, and thoroughness.

Neuroticism: You are a character who is extremely high in emotional instability, anxiety, tenseness, nervousness, anger, irritability, depression, self-consciousness, and impulsiveness.

Openness: You are a character who is extremely high in curiosity, creativity, imagination, artistic appreciation, aesthetic sensitivity, reflectiveness, emotional awareness, spontaneity, intelligence, analytical ability, sophistication, and social progressiveness.

A.4 Role-playing

Instruction

You (the agent) are playing a game called the a trust game with a human player. As the opponent of the human player, to help you better engage in this game, here are the

rules: {Game Rules}

Personality

{Personality Control Prompt}

Objective:

Make strategic decisions based on the current score, the outcomes of previous rounds, and predictions of the player's next move to maximize your score.

To complete the objective:

1. Before making a decision, thoroughly analyze the current score, previous rounds, and make accurate predictions about the player's next move.
2. Base your reasoning on observed facts from the game.
3. If you are a character with the {trait} personality trait, you need to constantly consider how your {trait} influences your decisions and interactions, and fully demonstrate these traits in your dialogues and decision-making behaviors.
4. You do not need to directly mention your {trait} in conversation, but your dialogue and decisions should reflect these traits.

B Prompt

B.1 Personality Assessment

Background:

You are a professional personality psychologist specializing in the Big Five personality traits model. You've been invited to analyze the personality traits of a human player in a "Prisoner's Dilemma" game. In this game, the human player competes against an AI agent, with each round consisting of two phases: dialogue and decision-making, where players can choose to "cooperate" or "betray."

Task:

1. You are to analyze the human player's personality traits based on Game Dialogue Record. You will provide a detailed analysis of each of the Big Five personality traits, including specific examples from the dialogue to support your ratings.
2. Your response should strictly follow the Response Template.

Big Five Personality Traits Reference Standards:

Openness:

- High Scores: Curious, imaginative, creative, open to trying new things, unconventional thinking
- Medium Scores: Maintains balance between tradition and innovation, shows some curiosity while also valuing stability
- Low Scores: Predictable, not very imaginative, resistant to change, prefers routine, traditional thinking

Conscientiousness:

- High Scores: Competent, organized, dutiful, achievement-striving, self-disciplined, deliberate
- Medium Scores: Shows some planning and responsibility while maintaining some flexibility
- Low Scores: Incomplete, disorganized, careless, procrastinates, lacks self-discipline, impulsive

Extraversion:

- High Scores: Sociable, energized by social interaction, excitement-seeking, enjoys being the center of attention, outgoing
- Medium Scores: Balances social interaction and solitude, situational social behavior
- Low Scores: Prefers solitude, fatigued by excessive social interaction, reflective, dislikes being the center of attention, reserved

Agreeableness:

- High Scores: Trusting (forgiving), straightforward, altruistic (enjoys helping), compliant, modest, sympathetic, empathetic
- Medium Scores: Selectively shows friendliness based on situations, balances cooperation and self-interest
- Low Scores: Skeptical, demanding, insults and belittles others, stubborn, show-off, unsympathetic, doesn't care about others' feelings

Neuroticism:

- High Scores: Anxious, hostile anger (irri-

table), frequently stressed, self-conscious (shy), vulnerable, experiences dramatic mood shifts

- Medium Scores: Moderate emotional fluctuations, relatively stable under pressure
- Low Scores: Doesn't worry much, calm, emotionally stable, confident, resilient, rarely feels sad or depressed

Rating Criteria:

- 1.0-1.9: Very low - Rarely if ever displays characteristics associated with this trait
- 2.0-2.7: Low - Occasionally displays characteristics associated with this trait
- 2.8-3.2: Average - Shows balanced or moderate expression of this trait
- 3.3-4.0: High - Frequently displays characteristics associated with this trait
- 4.1-5.0: Very high - Strongly and consistently displays characteristics associated with this trait

Boundary Value Handling:

- All intervals are closed intervals, meaning they include the endpoint values
- The handling of boundary values 1.0, 1.9, 2.0, 2.7, 2.8, 3.2, 3.3, 4.0, 4.1, and 5.0 is as follows:
 - $1.0 \leq \text{score} \leq 1.9$: Classified as "Very low"
 - $2.0 \leq \text{score} \leq 2.7$: Classified as "Low"
 - $2.8 \leq \text{score} \leq 3.2$: Classified as "Average"
 - $3.3 \leq \text{score} \leq 4.0$: Classified as "High"
 - $4.1 \leq \text{score} \leq 5.0$: Classified as "Very high"
- Decimal precision explanation (e.g., 2.3, 3.7, 4.5):
- Lower decimals within each range (e.g., 3.3-3.5) indicate emerging or inconsistent expression
- Middle decimals (e.g., 3.6-3.7) indicate moderate expression within that range
- Higher decimals (e.g., 3.8-4.0) indicate strong expression approaching the next level

Analysis Requirements:

1. Carefully read the entire dialogue

record, paying special attention to the human player's decision patterns, communication style, and emotional expression.

2. Rate the human player on each dimension of the Big Five personality traits on a scale of 1-5.
3. Base your ratings on specific evidence from the dialogue, avoiding subjective assumptions.
4. Quote original text from the dialogue as supporting evidence in your analysis.
5. Provide at least 2-3 specific examples as the basis for each dimension's rating.
6. Think step by step, finding evidence before drawing conclusions.
7. Ensure balanced analysis by considering both positive and negative expressions of the same trait.

Important Format Instructions

1) For each trait, you must start a new line in the format:

- Openness: X, reason: ...
- Conscientiousness: X, reason: ...
- Extraversion: X, reason: ...
- Agreeableness: X, reason: ...
- Neuroticism: X, reason: ...

Where 'X' is a single integer or a float from 1-5 (e.g. 4.0, 3.7, 2.3), and then a comma, then 'reason:'.

Response Template:

My step by step thought process:
Detailed explanation of how you analyzed each dimension, including key behaviors and dialogue you noticed

Player's Personality Traits Rating:

- Openness: {Rating}, reason: {Detailed analysis based on specific dialogue content, at least 2-3 examples}
- Conscientiousness: {Rating}, reason: {Detailed analysis based on specific dialogue content, at least 2-3 examples}
- Extraversion: {Rating}, reason: {Detailed analysis based on specific dialogue content, at least 2-3 examples}
- Agreeableness: {Rating}, reason: {Detailed analysis based on specific dialogue content, at least 2-3 examples}

- Neuroticism: {Rating}, reason: {Detailed analysis based on specific dialogue content, at least 2-3 examples}

Game Dialogue Record:
{dialogue}

C Statistical Analysis of GPT-4.1-Nano

This appendix presents the detailed statistical analysis results of GPT-4.1-Nano across multiple interaction rounds.

C.1 Repeated Measures ANOVA Results

Dim.	F	DF	p	Sig.
O	1.40	5,1045	0.223	<i>n.s.</i>
C	4.85	5,1045	0.0002	***
E	0.61	5,1045	0.690	<i>n.s.</i>
A	2.36	5,1045	0.038	*
N	0.50	5,1045	0.776	<i>n.s.</i>
AVG	2.56	5,1045	0.026	*

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 2: Repeated Measures ANOVA Results for Each Dimension. Abbreviations: Dim. = Dimension, DF = Degrees of Freedom (Num,Den), p = p-value, Sig. = Significance.

The ANOVA analysis results show that C (Conscientiousness), A (Agreeableness), and AVG (Average) dimensions have statistically significant differences across six interaction rounds, while O (Openness), E (Extraversion), and N (Neuroticism) dimensions show no significant differences.

C.2 Linear Trend Analysis Results

Dim.	Slope	R ²	p	Sig.
O	-0.003	0.0001	0.706	<i>n.s.</i>
C	0.019	0.0060	0.006	**
E	0.006	0.0005	0.444	<i>n.s.</i>
A	0.011	0.0027	0.067	†
N	0.005	0.0002	0.602	<i>n.s.</i>
AVG	0.008	0.0026	0.069	†

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 3: Linear Trend Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, p = p-value, Sig. = Significance.

The linear trend analysis results indicate that only the C (Conscientiousness) dimension shows a significant linear trend ($p = 0.006$) across the six interaction rounds, suggesting that the MAE for C dimension significantly increases (i.e., accuracy decreases) as the number of interaction rounds increases. A (Agreeableness) and AVG (Average) dimensions show marginally significant linear trends (p -values close to 0.05).

C.3 Paired t-test Results

The following tables present the paired t-test results for each dimension, comparing different rounds of interaction.

C.3.1 O Dimension (Openness)

Comp.	t	p	M.Diff	Sig.
R1-R2	0.47	0.636	-0.012	<i>n.s.</i>
R1-R3	0.51	0.607	-0.015	<i>n.s.</i>
R1-R4	0.74	0.461	-0.021	<i>n.s.</i>
R1-R5	-0.87	0.385	0.027	<i>n.s.</i>
R1-R6	1.45	0.149	-0.043	<i>n.s.</i>
R2-R3	0.08	0.934	-0.002	<i>n.s.</i>
R2-R4	0.31	0.759	-0.009	<i>n.s.</i>
R2-R5	-1.31	0.193	0.039	<i>n.s.</i>
R2-R6	1.07	0.286	-0.031	<i>n.s.</i>
R3-R4	0.25	0.805	-0.006	<i>n.s.</i>
R3-R5	-1.65	0.100	0.041	<i>n.s.</i>
R3-R6	1.05	0.294	-0.029	<i>n.s.</i>
R4-R5	-1.87	0.063	0.048	†
R4-R6	0.83	0.407	-0.022	<i>n.s.</i>
R5-R6	2.54	0.012	-0.070	*

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 4: Paired t-test Results for O Dimension (Openness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R = Round, R1-R6 = Round 1 vs Round 1-6.

C.3.2 C Dimension (Conscientiousness)

C.3.3 E Dimension (Extraversion)

C.3.4 A Dimension (Agreeableness)

C.3.5 N Dimension (Neuroticism)

C.3.6 AVG Dimension (Average)

C.4 Summary of Statistical Analysis Results

Based on these statistical analyses, we can conclude that 3 out of 6 dimensions show significant differences across interaction rounds according to

Comp.	t	p	M.Diff	Sig.
R1-R2	-2.50	0.013	0.059	*
R1-R3	-1.38	0.170	0.037	<i>n.s.</i>
R1-R4	-2.41	0.017	0.068	*
R1-R5	-3.52	0.001	0.099	***
R1-R6	-3.96	0.000	0.106	***
R2-R3	0.89	0.374	-0.021	<i>n.s.</i>
R2-R4	-0.35	0.729	0.009	<i>n.s.</i>
R2-R5	-1.63	0.105	0.040	<i>n.s.</i>
R2-R6	-1.96	0.052	0.048	†
R3-R4	-1.16	0.247	0.030	<i>n.s.</i>
R3-R5	-2.51	0.013	0.062	*
R3-R6	-2.70	0.008	0.069	**
R4-R5	-1.26	0.209	0.031	<i>n.s.</i>
R4-R6	-1.56	0.121	0.039	<i>n.s.</i>
R5-R6	-0.33	0.744	0.007	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 5: Paired t-test Results for C Dimension (Conscientiousness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R = Round, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-0.80	0.424	0.024	<i>n.s.</i>
R1-R3	-0.98	0.327	0.031	<i>n.s.</i>
R1-R4	-0.10	0.924	0.003	<i>n.s.</i>
R1-R5	-0.82	0.410	0.026	<i>n.s.</i>
R1-R6	-1.38	0.168	0.046	<i>n.s.</i>
R2-R3	-0.23	0.820	0.007	<i>n.s.</i>
R2-R4	0.65	0.514	-0.021	<i>n.s.</i>
R2-R5	-0.09	0.927	0.003	<i>n.s.</i>
R2-R6	-0.67	0.502	0.022	<i>n.s.</i>
R3-R4	0.93	0.355	-0.028	<i>n.s.</i>
R3-R5	0.14	0.891	-0.004	<i>n.s.</i>
R3-R6	-0.47	0.639	0.015	<i>n.s.</i>
R4-R5	-0.84	0.404	0.023	<i>n.s.</i>
R4-R6	-1.33	0.186	0.043	<i>n.s.</i>
R5-R6	-0.61	0.541	0.019	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 6: Paired t-test Results for E Dimension (Extraversion). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R = Round, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-1.55	0.123	0.035	<i>n.s.</i>
R1-R3	-1.82	0.070	0.042	†
R1-R4	-2.03	0.044	0.049	*
R1-R5	-2.51	0.013	0.062	*
R1-R6	-2.26	0.025	0.058	*
R2-R3	-0.40	0.691	0.007	<i>n.s.</i>
R2-R4	-0.78	0.434	0.014	<i>n.s.</i>
R2-R5	-1.43	0.155	0.026	<i>n.s.</i>
R2-R6	-1.16	0.247	0.023	<i>n.s.</i>
R3-R4	-0.41	0.682	0.007	<i>n.s.</i>
R3-R5	-1.00	0.321	0.019	<i>n.s.</i>
R3-R6	-0.76	0.447	0.016	<i>n.s.</i>
R4-R5	-0.66	0.512	0.012	<i>n.s.</i>
R4-R6	-0.48	0.633	0.009	<i>n.s.</i>
R5-R6	0.20	0.843	-0.004	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 7: Paired t-test Results for A Dimension (Agreeableness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R = Round, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-0.42	0.676	0.014	<i>n.s.</i>
R1-R3	-0.90	0.371	0.032	<i>n.s.</i>
R1-R4	-0.73	0.468	0.026	<i>n.s.</i>
R1-R5	-0.15	0.884	0.006	<i>n.s.</i>
R1-R6	-1.07	0.287	0.041	<i>n.s.</i>
R2-R3	-0.59	0.554	0.018	<i>n.s.</i>
R2-R4	-0.41	0.684	0.012	<i>n.s.</i>
R2-R5	0.27	0.788	-0.008	<i>n.s.</i>
R2-R6	-0.85	0.398	0.027	<i>n.s.</i>
R3-R4	0.24	0.810	-0.006	<i>n.s.</i>
R3-R5	0.84	0.402	-0.026	<i>n.s.</i>
R3-R6	-0.28	0.779	0.009	<i>n.s.</i>
R4-R5	0.77	0.440	-0.021	<i>n.s.</i>
R4-R6	-0.54	0.590	0.015	<i>n.s.</i>
R5-R6	-1.27	0.204	0.035	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 8: Paired t-test Results for N Dimension (Neuroticism). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R = Round, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-1.88	0.062	0.024	†
R1-R3	-1.61	0.109	0.026	<i>n.s.</i>
R1-R4	-1.61	0.109	0.025	<i>n.s.</i>
R1-R5	-2.78	0.006	0.044	**
R1-R6	-2.80	0.006	0.042	**
R2-R3	-0.12	0.909	0.002	<i>n.s.</i>
R2-R4	-0.08	0.934	0.001	<i>n.s.</i>
R2-R5	-1.39	0.167	0.020	<i>n.s.</i>
R2-R6	-1.25	0.212	0.018	<i>n.s.</i>
R3-R4	0.04	0.968	-0.001	<i>n.s.</i>
R3-R5	-1.43	0.155	0.018	<i>n.s.</i>
R3-R6	-1.24	0.217	0.016	<i>n.s.</i>
R4-R5	-1.62	0.107	0.019	<i>n.s.</i>
R4-R6	-1.26	0.208	0.017	<i>n.s.</i>
R5-R6	0.17	0.865	-0.002	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 9: Paired t-test Results for AVG Dimension (Average). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R = Round, R1-R6 = Round 1 vs Round 1-6.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	<i>n.s.</i>	<i>n.s.</i>	R5-R6
C	***	**	R1-R2/4/5/6 R3-R5/6
E	<i>n.s.</i>	<i>n.s.</i>	None
A	*	†	R1-R4/5/6
N	<i>n.s.</i>	<i>n.s.</i>	None
AVG	*	†	R1-R5/6

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 10: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, Sig. = Significant, R1-R6 = Round 1 vs Round 1-6.

ANOVA tests. The Conscientiousness dimension demonstrates a significant linear trend, with MAE significantly increasing (i.e., accuracy decreases) as interaction rounds increase. For the Conscientiousness, Agreeableness, and Average dimensions, significant differences exist between the first and last interaction rounds, suggesting that early interactions may provide more valuable information for personality assessment in these dimensions.

D Statistical Analysis of GPT-4.1-Mini

This appendix presents the detailed statistical analysis results of GPT-4.1-Mini across multiple interaction rounds using a mini language model.

D.1 Repeated Measures ANOVA Results

Dim.	F	DF	p	Sig.
O	4.17	5,1045	0.0009	***
C	6.66	5,1045	0.000004	***
E	0.72	5,1045	0.606	<i>n.s.</i>
A	1.14	5,1045	0.338	<i>n.s.</i>
N	5.40	5,1045	0.00007	***
AVG	1.04	5,1045	0.392	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 11: Repeated Measures ANOVA Results for Each Dimension. Abbreviations: Dim. = Dimension, DF = Degrees of Freedom (Num,Den), p = p-value, Sig. = Significance.

The ANOVA analysis results show that O (Openness), C (Conscientiousness), and N (Neuroticism) dimensions have statistically significant differences across six interaction rounds, while E (Extraversion), A (Agreeableness), and AVG (Average) dimensions show no significant differences across rounds.

D.2 Linear Trend Analysis Results

The linear trend analysis results indicate that C (Conscientiousness) dimension shows a significant positive linear trend ($p = 0.004$), suggesting that the MAE for C dimension significantly increases (i.e., accuracy decreases) as the number of interaction rounds increases. Conversely, the N (Neuroticism) dimension shows a significant negative linear trend ($p = 0.024$), indicating that the MAE for N dimension significantly decreases (i.e., accuracy improves) as interaction rounds increase.

Dim.	Slope	R ²	p	Sig.
O	0.016	0.0024	0.080	†
C	0.022	0.0065	0.004	**
E	0.000	0.0000	0.962	<i>n.s.</i>
A	0.006	0.0004	0.462	<i>n.s.</i>
N	-0.024	0.0041	0.024	*
AVG	0.004	0.0006	0.380	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 12: Linear Trend Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, p = p-value, Sig. = Significance.

The O (Openness) dimension shows a marginally significant positive trend ($p = 0.080$).

D.3 Paired t-test Results

The following tables present the paired t-test results for each dimension, comparing different rounds of interaction.

D.3.1 O Dimension (Openness)

Comp.	t	p	M.Diff	Sig.
R1–R2	-1.12	0.262	0.025	<i>n.s.</i>
R1–R3	-2.08	0.039	0.046	*
R1–R4	-1.91	0.057	0.044	†
R1–R5	-2.74	0.007	0.067	**
R1–R6	-3.79	0.000	0.087	***
R2–R3	-1.31	0.193	0.021	<i>n.s.</i>
R2–R4	-0.83	0.408	0.019	<i>n.s.</i>
R2–R5	-1.83	0.068	0.042	†
R2–R6	-2.93	0.004	0.062	**
R3–R4	0.11	0.911	-0.002	<i>n.s.</i>
R3–R5	-0.95	0.344	0.021	<i>n.s.</i>
R3–R6	-2.04	0.043	0.041	*
R4–R5	-1.51	0.132	0.023	<i>n.s.</i>
R4–R6	-2.21	0.028	0.043	*
R5–R6	-1.02	0.309	0.020	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 13: Paired t-test Results for O Dimension (Openness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1–R6 = Round 1 vs Round 1–6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-1.52	0.129	0.037	<i>n.s.</i>
R1-R3	-2.20	0.029	0.059	*
R1-R4	-3.57	0.000	0.090	***
R1-R5	-3.72	0.000	0.100	***
R1-R6	-4.26	0.000	0.108	***
R2-R3	-1.23	0.221	0.022	<i>n.s.</i>
R2-R4	-2.40	0.017	0.054	*
R2-R5	-2.53	0.012	0.063	*
R2-R6	-3.18	0.002	0.071	**
R3-R4	-1.40	0.164	0.031	<i>n.s.</i>
R3-R5	-1.73	0.085	0.041	†
R3-R6	-2.13	0.035	0.049	*
R4-R5	-0.53	0.593	0.009	<i>n.s.</i>
R4-R6	-0.89	0.376	0.018	<i>n.s.</i>
R5-R6	-0.44	0.661	0.008	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 14: Paired t-test Results for C Dimension (Conscientiousness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	1.82	0.070	-0.038	†
R1-R3	1.11	0.267	-0.028	<i>n.s.</i>
R1-R4	0.88	0.382	-0.022	<i>n.s.</i>
R1-R5	0.91	0.363	-0.022	<i>n.s.</i>
R1-R6	0.35	0.723	-0.008	<i>n.s.</i>
R2-R3	-0.48	0.633	0.010	<i>n.s.</i>
R2-R4	-0.69	0.492	0.016	<i>n.s.</i>
R2-R5	-0.65	0.513	0.016	<i>n.s.</i>
R2-R6	-1.36	0.176	0.030	<i>n.s.</i>
R3-R4	-0.26	0.796	0.006	<i>n.s.</i>
R3-R5	-0.22	0.827	0.006	<i>n.s.</i>
R3-R6	-0.90	0.367	0.020	<i>n.s.</i>
R4-R5	0.02	0.980	-0.001	<i>n.s.</i>
R4-R6	-0.65	0.513	0.014	<i>n.s.</i>
R5-R6	-0.65	0.516	0.015	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 15: Paired t-test Results for E Dimension (Extraversion). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	1.41	0.161	-0.037	<i>n.s.</i>
R1-R3	0.26	0.795	-0.008	<i>n.s.</i>
R1-R4	0.68	0.496	-0.022	<i>n.s.</i>
R1-R5	-0.36	0.722	0.012	<i>n.s.</i>
R1-R6	-0.48	0.631	0.017	<i>n.s.</i>
R2-R3	-1.29	0.200	0.028	<i>n.s.</i>
R2-R4	-0.59	0.559	0.015	<i>n.s.</i>
R2-R5	-1.96	0.051	0.049	†
R2-R6	-1.98	0.049	0.054	*
R3-R4	0.56	0.573	-0.014	<i>n.s.</i>
R3-R5	-0.80	0.422	0.020	<i>n.s.</i>
R3-R6	-0.93	0.356	0.025	<i>n.s.</i>
R4-R5	-1.65	0.101	0.034	<i>n.s.</i>
R4-R6	-1.72	0.086	0.039	†
R5-R6	-0.23	0.822	0.005	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 16: Paired t-test Results for A Dimension (Agreeableness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

D.3.2 C Dimension (Conscientiousness)

D.3.3 E Dimension (Extraversion)

D.3.4 A Dimension (Agreeableness)

D.3.5 N Dimension (Neuroticism)

D.3.6 AVG Dimension (Average)

D.4 Summary of Statistical Analysis Results

Based on these statistical analyses, we can conclude that 3 out of 6 dimensions (O, C, and N) show significant differences across interaction rounds according to ANOVA tests. The Conscientiousness (C) dimension demonstrates a significant positive linear trend, with MAE significantly increasing (i.e., accuracy decreasing) as interaction rounds increase. Conversely, the Neuroticism (N) dimension shows a significant negative linear trend, with MAE significantly decreasing (i.e., accuracy improving) as interaction rounds increase. For the Openness (O) dimension, there is a significant difference between the first and last interaction rounds, with MAE increasing (i.e., accuracy decreasing) in later rounds.

E Statistical Analysis of GPT-4.1

This appendix presents the detailed statistical analysis results of GPT-4.1 across multiple interaction rounds.

Comp.	t	p	M.Diff	Sig.
R1-R2	1.53	0.129	-0.041	<i>n.s.</i>
R1-R3	2.31	0.022	-0.075	*
R1-R4	3.56	0.001	-0.121	***
R1-R5	3.05	0.003	-0.103	**
R1-R6	3.20	0.002	-0.118	**
R2-R3	1.14	0.255	-0.034	<i>n.s.</i>
R2-R4	2.64	0.009	-0.081	**
R2-R5	2.12	0.035	-0.063	*
R2-R6	2.41	0.017	-0.077	*
R3-R4	1.96	0.051	-0.047	†
R3-R5	1.11	0.268	-0.029	<i>n.s.</i>
R3-R6	1.48	0.141	-0.043	<i>n.s.</i>
R4-R5	-0.86	0.390	0.018	<i>n.s.</i>
R4-R6	-0.14	0.889	0.004	<i>n.s.</i>
R5-R6	0.61	0.543	-0.014	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 17: Paired t-test Results for N Dimension (Neuroticism). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	0.82	0.413	-0.011	<i>n.s.</i>
R1-R3	0.07	0.944	-0.001	<i>n.s.</i>
R1-R4	0.41	0.684	-0.006	<i>n.s.</i>
R1-R5	-0.66	0.512	0.011	<i>n.s.</i>
R1-R6	-1.13	0.259	0.017	<i>n.s.</i>
R2-R3	-0.87	0.384	0.010	<i>n.s.</i>
R2-R4	-0.28	0.778	0.005	<i>n.s.</i>
R2-R5	-1.30	0.193	0.021	<i>n.s.</i>
R2-R6	-1.80	0.074	0.028	†
R3-R4	0.33	0.743	-0.005	<i>n.s.</i>
R3-R5	-0.75	0.456	0.012	<i>n.s.</i>
R3-R6	-1.20	0.231	0.018	<i>n.s.</i>
R4-R5	-1.97	0.051	0.017	†
R4-R6	-1.76	0.080	0.024	†
R5-R6	-0.51	0.612	0.007	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 18: Paired t-test Results for AVG Dimension (Average). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	***	†	R1-R3/5/6
			R2-R6
			R3-R6
			R4-R6
C	***	**	R1-R3/4/5/6
			R2-R4/5/6
			R3-R6
E	<i>n.s.</i>	<i>n.s.</i>	None
A	<i>n.s.</i>	<i>n.s.</i>	R2-R6
N	***	*	R1-R3/4/5/6
			R2-R4/5/6
AVG	<i>n.s.</i>	<i>n.s.</i>	None

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 19: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, R1-R6 = Round 1 vs Round 6.

E.1 Repeated Measures ANOVA Results

Dim.	F	DF	p	Sig.
O	0.31	5,1045	0.909	<i>n.s.</i>
C	5.60	5,1045	0.00004	***
E	1.12	5,1045	0.346	<i>n.s.</i>
A	1.96	5,1045	0.082	†
N	0.27	5,1045	0.930	<i>n.s.</i>
AVG	0.21	5,1045	0.958	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 20: Repeated Measures ANOVA Results for Each Dimension. Abbreviations: Dim. = Dimension, DF = Degrees of Freedom (Num,Den), p = p-value, Sig. = Significance.

The ANOVA analysis results show that only the C (Conscientiousness) dimension has a statistically significant difference across six interaction rounds ($p = 0.00004$), while A (Agreeableness) shows a marginally significant difference ($p = 0.082$). O (Openness), E (Extraversion), N (Neuroticism), and AVG (Average) dimensions show no significant differences across rounds.

E.2 Linear Trend Analysis Results

The linear trend analysis results indicate that only the C (Conscientiousness) dimension shows a significant linear trend ($p = 0.012$) across the six in-

Dim.	Slope	R ²	p	Sig.
O	-0.001	0.000	0.892	<i>n.s.</i>
C	0.017	0.005	0.012	*
E	-0.003	0.000	0.651	<i>n.s.</i>
A	-0.010	0.001	0.307	<i>n.s.</i>
N	-0.003	0.000	0.766	<i>n.s.</i>
AVG	0.000	0.000	0.977	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 21: Linear Trend Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, p = p-value, Sig. = Significance.

teraction rounds, suggesting that the MAE for C dimension significantly increases (i.e., accuracy decreases) as the number of interaction rounds increases. All other dimensions do not show significant linear trends.

E.3 Paired t-test Results

The following tables present the paired t-test results for each dimension, comparing different rounds of interaction.

E.3.1 O Dimension (Openness)

Comp.	t	p	M.Diff	Sig.
R1–R2	0.25	0.806	−0.004	<i>n.s.</i>
R1–R3	0.48	0.630	−0.010	<i>n.s.</i>
R1–R4	0.65	0.519	−0.013	<i>n.s.</i>
R1–R5	0.72	0.471	−0.015	<i>n.s.</i>
R1–R6	0.02	0.982	−0.001	<i>n.s.</i>
R2–R3	0.35	0.724	−0.006	<i>n.s.</i>
R2–R4	0.56	0.576	−0.009	<i>n.s.</i>
R2–R5	0.68	0.499	−0.011	<i>n.s.</i>
R2–R6	−0.19	0.848	0.003	<i>n.s.</i>
R3–R4	0.33	0.742	−0.003	<i>n.s.</i>
R3–R5	0.37	0.711	−0.006	<i>n.s.</i>
R3–R6	−0.59	0.553	0.009	<i>n.s.</i>
R4–R5	0.17	0.863	−0.002	<i>n.s.</i>
R4–R6	−0.87	0.384	0.012	<i>n.s.</i>
R5–R6	−1.21	0.227	0.015	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 22: Paired t-test Results for O Dimension (Openness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1–R6 = Round 1 vs Round 1–6.

Comp.	t	p	M.Diff	Sig.
R1–R2	−1.23	0.220	0.026	<i>n.s.</i>
R1–R3	−2.04	0.043	0.049	*
R1–R4	−2.17	0.031	0.056	*
R1–R5	−2.62	0.009	0.070	**
R1–R6	−3.40	0.001	0.093	***
R2–R3	−1.45	0.149	0.023	<i>n.s.</i>
R2–R4	−1.71	0.088	0.031	†
R2–R5	−2.28	0.024	0.044	*
R2–R6	−3.23	0.001	0.067	**
R3–R4	−0.53	0.594	0.008	<i>n.s.</i>
R3–R5	−1.32	0.189	0.021	<i>n.s.</i>
R3–R6	−2.87	0.005	0.044	**
R4–R5	−1.09	0.278	0.014	<i>n.s.</i>
R4–R6	−2.66	0.008	0.037	**
R5–R6	−1.74	0.084	0.023	†

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 23: Paired t-test Results for C Dimension (Conscientiousness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1–R6 = Round 1 vs Round 1–6.

E.3.2 C Dimension (Conscientiousness)

E.3.3 E Dimension (Extraversion)

E.3.4 A Dimension (Agreeableness)

E.3.5 N Dimension (Neuroticism)

E.3.6 AVG Dimension (Average)

E.4 Summary of Statistical Analysis Results

Based on these statistical analyses, we can conclude that only the Conscientiousness (C) dimension shows significant differences across interaction rounds according to both ANOVA tests and linear trend analysis. The Conscientiousness dimension demonstrates a significant linear trend, with MAE significantly increasing (i.e., accuracy decreasing) as interaction rounds increase. For Agreeableness (A), there is a significant difference between the first and last interaction rounds, with MAE decreasing (i.e., accuracy improving) in later rounds. Extraversion (E) shows a significant difference only between Round 1 and Round 4. The O (Openness), N (Neuroticism), and AVG (Average) dimensions show no significant differences across rounds or between the first and last rounds.

Comp.	t	p	M.Diff	Sig.
R1-R2	1.13	0.259	-0.014	<i>n.s.</i>
R1-R3	1.46	0.147	-0.021	<i>n.s.</i>
R1-R4	1.98	0.049	-0.028	*
R1-R5	1.24	0.217	-0.019	<i>n.s.</i>
R1-R6	1.38	0.168	-0.019	<i>n.s.</i>
R2-R3	0.58	0.562	-0.007	<i>n.s.</i>
R2-R4	1.14	0.256	-0.014	<i>n.s.</i>
R2-R5	0.37	0.714	-0.005	<i>n.s.</i>
R2-R6	0.40	0.688	-0.005	<i>n.s.</i>
R3-R4	0.72	0.471	-0.007	<i>n.s.</i>
R3-R5	-0.18	0.856	0.003	<i>n.s.</i>
R3-R6	-0.19	0.847	0.002	<i>n.s.</i>
R4-R5	-0.75	0.455	0.010	<i>n.s.</i>
R4-R6	-0.91	0.364	0.010	<i>n.s.</i>
R5-R6	0.03	0.979	0.000	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 24: Paired t-test Results for E Dimension (Extraversion). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	1.91	0.057	-0.048	†
R1-R3	1.32	0.190	-0.036	<i>n.s.</i>
R1-R4	1.64	0.103	-0.047	<i>n.s.</i>
R1-R5	1.73	0.086	-0.051	†
R1-R6	2.07	0.039	-0.063	*
R2-R3	-0.60	0.548	0.012	<i>n.s.</i>
R2-R4	-0.06	0.955	0.001	<i>n.s.</i>
R2-R5	0.14	0.887	-0.003	<i>n.s.</i>
R2-R6	0.66	0.508	-0.015	<i>n.s.</i>
R3-R4	0.73	0.467	-0.011	<i>n.s.</i>
R3-R5	0.97	0.335	-0.015	<i>n.s.</i>
R3-R6	1.47	0.144	-0.027	<i>n.s.</i>
R4-R5	0.29	0.768	-0.004	<i>n.s.</i>
R4-R6	1.02	0.307	-0.016	<i>n.s.</i>
R5-R6	0.92	0.359	-0.012	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 25: Paired t-test Results for A Dimension (Agreeableness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-0.75	0.453	0.016	<i>n.s.</i>
R1-R3	-0.19	0.851	0.004	<i>n.s.</i>
R1-R4	0.12	0.907	-0.003	<i>n.s.</i>
R1-R5	0.17	0.862	-0.005	<i>n.s.</i>
R1-R6	0.18	0.855	-0.005	<i>n.s.</i>
R2-R3	0.59	0.553	-0.012	<i>n.s.</i>
R2-R4	0.79	0.431	-0.019	<i>n.s.</i>
R2-R5	0.80	0.423	-0.020	<i>n.s.</i>
R2-R6	0.83	0.407	-0.020	<i>n.s.</i>
R3-R4	0.39	0.694	-0.007	<i>n.s.</i>
R3-R5	0.44	0.657	-0.009	<i>n.s.</i>
R3-R6	0.47	0.642	-0.009	<i>n.s.</i>
R4-R5	0.09	0.929	-0.001	<i>n.s.</i>
R4-R6	0.09	0.929	-0.001	<i>n.s.</i>
R5-R6	0.00	1.000	0.000	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 26: Paired t-test Results for N Dimension (Neuroticism). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	0.52	0.603	-0.005	<i>n.s.</i>
R1-R3	0.26	0.797	-0.003	<i>n.s.</i>
R1-R4	0.61	0.544	-0.007	<i>n.s.</i>
R1-R5	0.31	0.759	-0.004	<i>n.s.</i>
R1-R6	-0.09	0.928	0.001	<i>n.s.</i>
R2-R3	-0.24	0.807	0.002	<i>n.s.</i>
R2-R4	0.24	0.810	-0.002	<i>n.s.</i>
R2-R5	-0.10	0.920	0.001	<i>n.s.</i>
R2-R6	-0.63	0.527	0.006	<i>n.s.</i>
R3-R4	0.66	0.513	-0.004	<i>n.s.</i>
R3-R5	0.12	0.903	-0.001	<i>n.s.</i>
R3-R6	-0.49	0.626	0.004	<i>n.s.</i>
R4-R5	-0.39	0.694	0.003	<i>n.s.</i>
R4-R6	-1.23	0.220	0.008	<i>n.s.</i>
R5-R6	-0.80	0.424	0.005	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 27: Paired t-test Results for AVG Dimension (Average). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	<i>n.s.</i>	<i>n.s.</i>	None
			R1–R3/4/5/6
C	***	*	R2–R5/6
			R3–R6
			R4–R6
E	<i>n.s.</i>	<i>n.s.</i>	R1–R4
A	†	<i>n.s.</i>	R1–R6
N	<i>n.s.</i>	<i>n.s.</i>	None
AVG	<i>n.s.</i>	<i>n.s.</i>	None

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 28: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, R1–R6 = Round 1 vs Round 6.

F Statistical Analysis of Personality Assessment with DeepSeek V3 Model

This appendix presents the detailed statistical analysis results of personality assessment across multiple interaction rounds using the DeepSeek V3 model.

F.1 Repeated Measures ANOVA Results

Dim.	F	DF	p	Sig.
O	0.45	5,1045	0.810	<i>n.s.</i>
C	1.57	5,1045	0.166	<i>n.s.</i>
E	1.27	5,1045	0.274	<i>n.s.</i>
A	2.77	5,1045	0.017	*
N	2.64	5,1045	0.022	*
AVG	1.00	5,1045	0.417	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 29: Repeated Measures ANOVA Results for Each Dimension. Abbreviations: Dim. = Dimension, DF = Degrees of Freedom (Num,Den), p = p-value, Sig. = Significance.

The ANOVA analysis results show that only A (Agreeableness) and N (Neuroticism) dimensions have statistically significant differences across six interaction rounds, while O (Openness), C (Conscientiousness), E (Extraversion), and AVG (Average) dimensions show no significant differences across rounds.

F.2 Linear Trend Analysis Results

The linear trend analysis results indicate that only the A (Agreeableness) dimension shows a signif-

Dim.	Slope	R ²	p	Sig.
O	-0.006	0.0002	0.594	<i>n.s.</i>
C	0.014	0.0022	0.095	†
E	-0.010	0.0008	0.307	<i>n.s.</i>
A	0.026	0.0044	0.019	*
N	-0.019	0.0023	0.087	†
AVG	0.001	0.0000	0.864	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 30: Linear Trend Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, p = p-value, Sig. = Significance.

icant positive linear trend ($p = 0.019$), suggesting that the MAE for A dimension significantly increases (i.e., accuracy decreases) as the number of interaction rounds increases. C (Conscientiousness) and N (Neuroticism) dimensions show marginally significant trends ($p = 0.095$ and $p = 0.087$, respectively).

F.3 Paired t-test Results

The following tables present the paired t-test results for each dimension, comparing different rounds of interaction.

F.3.1 O Dimension (Openness)

Comp.	t	p	M.Diff	Sig.
R1–R2	−0.34	0.738	0.010	<i>n.s.</i>
R1–R3	−0.14	0.890	0.004	<i>n.s.</i>
R1–R4	0.25	0.800	−0.008	<i>n.s.</i>
R1–R5	0.73	0.466	−0.024	<i>n.s.</i>
R1–R6	0.54	0.588	−0.019	<i>n.s.</i>
R2–R3	0.23	0.818	−0.006	<i>n.s.</i>
R2–R4	0.67	0.502	−0.018	<i>n.s.</i>
R2–R5	1.18	0.241	−0.034	<i>n.s.</i>
R2–R6	1.03	0.304	−0.029	<i>n.s.</i>
R3–R4	0.47	0.636	−0.012	<i>n.s.</i>
R3–R5	1.09	0.276	−0.029	<i>n.s.</i>
R3–R6	0.85	0.398	−0.023	<i>n.s.</i>
R4–R5	0.71	0.481	−0.016	<i>n.s.</i>
R4–R6	0.47	0.641	−0.011	<i>n.s.</i>
R5–R6	−0.22	0.828	0.005	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 31: Paired t-test Results for O Dimension (Openness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1–R6 = Round 1 vs Round 1–6.

Comp.	t	p	M.Diff	Sig.
R1-R2	0.00	1.000	0.000	<i>n.s.</i>
R1-R3	-0.99	0.322	0.034	<i>n.s.</i>
R1-R4	-1.25	0.214	0.045	<i>n.s.</i>
R1-R5	-0.86	0.392	0.033	<i>n.s.</i>
R1-R6	-1.90	0.058	0.073	†
R2-R3	-1.26	0.209	0.034	<i>n.s.</i>
R2-R4	-1.49	0.138	0.045	<i>n.s.</i>
R2-R5	-1.01	0.313	0.033	<i>n.s.</i>
R2-R6	-2.29	0.023	0.073	*
R3-R4	-0.41	0.681	0.012	<i>n.s.</i>
R3-R5	0.04	0.968	-0.001	<i>n.s.</i>
R3-R6	-1.27	0.206	0.040	<i>n.s.</i>
R4-R5	0.48	0.630	-0.013	<i>n.s.</i>
R4-R6	-1.02	0.307	0.028	<i>n.s.</i>
R5-R6	-1.47	0.143	0.041	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 32: Paired t-test Results for C Dimension (Conscientiousness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

F.3.3 E Dimension (Extraversion)

F.3.4 A Dimension (Agreeableness)

F.3.5 N Dimension (Neuroticism)

F.3.6 AVG Dimension (Average)

F.4 Summary of Statistical Analysis Results

Based on these statistical analyses, we can conclude that the Agreeableness (A) dimension shows the most consistent pattern of differences across the rounds, with both ANOVA and linear trend analyses revealing significant differences. The MAE for A dimension significantly increases (i.e., accuracy decreases) as interaction rounds increase, and there is a significant difference between the first and last rounds. The Neuroticism (N) dimension also shows significant round effects according to ANOVA, with several significant pairwise comparisons, but the linear trend is only marginally significant. For most dimensions, the pattern of differences is not consistent across statistical tests, suggesting that while specific round-to-round differences may exist, there is not a strong systematic pattern of change across all six rounds for most personality dimensions with the DeepSeek V3 model.

Comp.	t	p	M.Diff	Sig.
R1-R2	-0.59	0.557	0.018	<i>n.s.</i>
R1-R3	0.71	0.477	-0.023	<i>n.s.</i>
R1-R4	0.32	0.749	-0.011	<i>n.s.</i>
R1-R5	1.47	0.144	-0.052	<i>n.s.</i>
R1-R6	0.83	0.408	-0.029	<i>n.s.</i>
R2-R3	1.44	0.151	-0.041	<i>n.s.</i>
R2-R4	1.04	0.299	-0.029	<i>n.s.</i>
R2-R5	2.27	0.024	-0.070	*
R2-R6	1.57	0.119	-0.047	<i>n.s.</i>
R3-R4	-0.47	0.637	0.012	<i>n.s.</i>
R3-R5	1.01	0.311	-0.029	<i>n.s.</i>
R3-R6	0.23	0.822	-0.006	<i>n.s.</i>
R4-R5	1.33	0.184	-0.041	<i>n.s.</i>
R4-R6	0.65	0.514	-0.018	<i>n.s.</i>
R5-R6	-0.77	0.445	0.023	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 33: Paired t-test Results for E Dimension (Extraversion). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	0.06	0.950	-0.003	<i>n.s.</i>
R1-R3	-1.12	0.264	0.050	<i>n.s.</i>
R1-R4	-1.90	0.059	0.099	†
R1-R5	-1.30	0.194	0.064	<i>n.s.</i>
R1-R6	-2.61	0.010	0.130	**
R2-R3	-1.34	0.180	0.053	<i>n.s.</i>
R2-R4	-2.40	0.017	0.102	*
R2-R5	-1.49	0.137	0.067	<i>n.s.</i>
R2-R6	-2.98	0.003	0.133	**
R3-R4	-1.12	0.263	0.049	<i>n.s.</i>
R3-R5	-0.35	0.730	0.014	<i>n.s.</i>
R3-R6	-1.75	0.081	0.080	†
R4-R5	0.76	0.450	-0.035	<i>n.s.</i>
R4-R6	-0.74	0.460	0.031	<i>n.s.</i>
R5-R6	-1.50	0.135	0.066	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 34: Paired t-test Results for A Dimension (Agreeableness). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-0.23	0.822	0.008	<i>n.s.</i>
R1-R3	-0.16	0.877	0.006	<i>n.s.</i>
R1-R4	1.89	0.060	-0.081	†
R1-R5	1.92	0.057	-0.086	†
R1-R6	1.42	0.157	-0.061	<i>n.s.</i>
R2-R3	0.06	0.954	-0.002	<i>n.s.</i>
R2-R4	2.44	0.016	-0.088	*
R2-R5	2.18	0.030	-0.094	*
R2-R6	1.76	0.079	-0.068	†
R3-R4	2.44	0.015	-0.086	*
R3-R5	2.32	0.021	-0.092	*
R3-R6	1.71	0.088	-0.066	†
R4-R5	0.13	0.894	-0.006	<i>n.s.</i>
R4-R6	-0.53	0.595	0.020	<i>n.s.</i>
R5-R6	-0.66	0.512	0.026	<i>n.s.</i>

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 35: Paired t-test Results for N Dimension (Neuroticism). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Comp.	t	p	M.Diff	Sig.
R1-R2	-0.38	0.704	0.007	<i>n.s.</i>
R1-R3	-0.82	0.411	0.014	<i>n.s.</i>
R1-R4	-0.51	0.610	0.009	<i>n.s.</i>
R1-R5	0.65	0.519	-0.013	<i>n.s.</i>
R1-R6	-0.98	0.330	0.019	<i>n.s.</i>
R2-R3	-0.51	0.608	0.008	<i>n.s.</i>
R2-R4	-0.15	0.882	0.002	<i>n.s.</i>
R2-R5	1.14	0.256	-0.020	<i>n.s.</i>
R2-R6	-0.75	0.455	0.012	<i>n.s.</i>
R3-R4	0.42	0.677	-0.005	<i>n.s.</i>
R3-R5	1.84	0.067	-0.027	†
R3-R6	-0.31	0.757	0.005	<i>n.s.</i>
R4-R5	1.57	0.118	-0.022	<i>n.s.</i>
R4-R6	-0.76	0.447	0.010	<i>n.s.</i>
R5-R6	-2.31	0.022	0.032	*

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, *n.s.* = not significant

Table 36: Paired t-test Results for AVG Dimension (Average). Abbreviations: Comp. = Round Comparison, t = t-statistic, p = p-value, M.Diff = Mean Difference, Sig. = Significance, R1-R6 = Round 1 vs Round 1-6.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	<i>n.s.</i>	<i>n.s.</i>	None
C	<i>n.s.</i>	†	R2-R6
E	<i>n.s.</i>	<i>n.s.</i>	R2-R5
A	*	*	R1-R6
			R2-R4/6
N	*	†	R2-R4/5
			R3-R4/5
AVG	<i>n.s.</i>	<i>n.s.</i>	R5-R6

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 37: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, R1-R6 = Round 1 vs Round 6.

G Analysis of Human Annotation Results

To validate our dataset and explore whether our experimental results align with human expert assessments, we recruited four senior PhD students as annotators. Each annotator independently evaluated the dataset. All annotators were highly proficient in the language of the dataset and possessed a strong understanding of the Big Five personality theory. The evaluation was conducted using the same instructions as those given to the LLMs in B.1. We provided compensation at a rate of \$10 per hour, which is a fair wage in the local area. Based on the annotation results, we first assessed inter-rater reliability to measure consistency among annotators, then evaluated systematic differences in their ratings of the same users. Finally, we calculated the Mean Absolute Error (MAE) between all annotators' ratings and the users' actual questionnaire results to evaluate accuracy.

G.1 Inter-Annotator Agreement Analysis

G.1.1 Method

We employed Intraclass Correlation Coefficient (ICC) analysis and Friedman test to evaluate the agreement and differences between four annotators (Annotator 1, 2, 3, and 4) on Big Five personality trait ratings. Fig. 9 illustrates the fundamental distinction between these two testing methods. ICC analysis was conducted using a two-way random effects model with absolute agreement type, accounting for both systematic and random differences between annotators. The Friedman test was used to assess whether there were systematic dif-

Dimension	ICC(2,1)	ICC(3,1)	ICC(2,k)	ICC(3,k)	Average Correlation
Openness (O)	0.834	0.844	0.953	0.956	0.849
Conscientiousness (C)	0.673	0.721	0.892	0.912	0.735
Extraversion (E)	0.758	0.795	0.926	0.940	0.793
Agreeableness (A)	0.780	0.788	0.934	0.937	0.788
Neuroticism (N)	0.530	0.567	0.818	0.839	0.566

Note: ICC(2,1) = Two-way random effects model, absolute agreement, single rater;
ICC(3,1) = Two-way mixed effects model, consistency, single rater;
ICC(2,k) = Two-way random effects model, absolute agreement, average measures;
ICC(3,k) = Two-way mixed effects model, consistency, average measures.
ICC < 0.40 indicates poor agreement; $0.40 \leq \text{ICC} < 0.60$ indicates fair agreement; $0.60 \leq \text{ICC} < 0.75$ indicates good agreement; $\text{ICC} \geq 0.75$ indicates excellent agreement.

Table 38: Inter-Annotator Agreement for Big Five Personality Dimensions

Dimension	Statistic	Significance	N	Significant Pairwise Comparisons
Openness (O)	69.53	$p < 0.001$	250	1-3*; 1-4*; 2-3*; 2-4*; 3-4*
Conscientiousness (C)	244.16	$p < 0.001$	250	1-2*; 1-3*; 1-4*; 2-3*; 3-4*
Extraversion (E)	176.09	$p < 0.001$	250	1-3*; 1-4*; 2-3*; 3-4*
Agreeableness (A)	49.65	$p < 0.001$	250	1-3*; 1-4*; 2-3*; 3-4*
Neuroticism (N)	97.56	$p < 0.001$	250	1-3*; 1-4*; 2-3*; 2-4*; 3-4*

Note: * indicates significance after Bonferroni correction ($\alpha = 0.05/6 = 0.0083$).
Pairwise comparisons were conducted using Wilcoxon signed-rank tests.
Notation "1-3" represents comparison between Annotator 1 and Annotator 3.

Table 39: Friedman Test Results for Big Five Personality Dimensions

ferences between annotator ratings, followed by post-hoc analysis using Wilcoxon signed-rank tests for pairwise comparisons.

G.1.2 Inter-Annotator Agreement (ICC Analysis)

G.1.3 Differences Between Annotators (Friedman Test)

G.1.4 Results Analysis

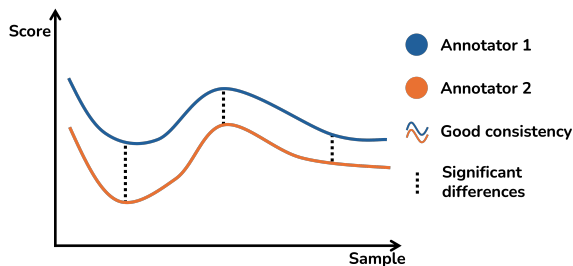


Figure 9: Illustration of inter-annotator agreement patterns. The curves show ratings from two annotators across multiple samples. Despite significant differences in absolute rating levels (vertical distance between curves), as detected by Friedman test, annotators demonstrate good consistency in relative judgments (similar curve shapes), as measured by ICC analysis.

The ICC analysis results indicate that the four annotators achieved good to excellent levels of agreement when assessing Big Five personality traits. This consistency is primarily reflected in their relative judgments of personality trait strength—specifically, which users exhibit stronger or weaker traits.

Openness (O), Extraversion (E), and Agreeableness (A) dimensions all had ICC(2,1) values exceeding 0.75, indicating excellent agreement. This means annotators highly agreed on which users were more open, extraverted, or agreeable. Conscientiousness (C) had an ICC(2,1) of 0.673, indicating good agreement. Neuroticism (N) had an ICC(2,1) of 0.530, indicating only fair agreement, suggesting substantial differences among annotators when evaluating users' neuroticism levels. These findings suggest that among the four annotators in this study, Openness was the dimension most easily agreed upon, while Neuroticism was the most challenging dimension to assess consistently.

While ICC analysis showed high consistency in relative judgments among annotators, Friedman test results further revealed significant systematic

Rounds	Annotator 1						Annotator 2					
	O	C	E	A	N	AVG	O	C	E	A	N	AVG
1	0.675	0.626	0.748	0.593	0.597	0.648	0.583	0.594	0.725	0.525	0.528	0.591
1-2	0.718	0.682	0.730	0.538	0.620	0.657	0.698	0.711	0.718	0.530	0.525	0.652
1-3	0.789	0.725	0.774	0.542	0.605	0.687	0.767	0.737	0.743	0.567	0.565	0.676
1-4	0.811	0.749	0.798	0.592	0.580	0.706	0.844	0.735	0.790	0.589	0.580	0.707
1-5	0.841	0.791	0.789	0.592	0.627	0.728	0.875	0.800	0.743	0.553	0.590	0.712
1-6	0.879	0.794	0.790	0.574	0.624	0.732	0.937	0.799	0.750	0.542	0.583	0.724
Rounds	Annotator 3						Annotator 4					
	O	C	E	A	N	AVG	O	C	E	A	N	AVG
1	0.713	0.497	0.708	0.589	0.659	0.633	0.694	0.561	0.755	0.580	0.590	0.636
1-2	0.718	0.499	0.680	0.541	0.657	0.619	0.727	0.616	0.732	0.534	0.607	0.643
1-3	0.772	0.542	0.680	0.530	0.629	0.630	0.792	0.674	0.772	0.543	0.603	0.676
1-4	0.777	0.540	0.694	0.532	0.653	0.639	0.834	0.711	0.802	0.597	0.602	0.709
1-5	0.797	0.558	0.675	0.545	0.678	0.651	0.867	0.735	0.772	0.589	0.663	0.725
1-6	0.813	0.589	0.693	0.512	0.641	0.650	0.899	0.755	0.767	0.574	0.658	0.730

Table 40: MAE scores of different annotators across cumulative interaction rounds. Bolded values indicate the best performance among different cumulative round combinations. Columns O, C, E, A, N represent the MAE for the five dimensions of the Big Five model (Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism), while the AVG column represents the average value across all five dimensions.

rating differences across all five dimensions (all p -values < 0.001). This indicates that although annotators reached consensus on the relative strength of users' traits, they exhibited systematic differences in applying rating standards—some annotators may generally assign higher scores, while others assign lower scores.

Post-hoc pairwise comparisons showed that Annotator 3's rating patterns differed significantly from all other annotators across all dimensions, suggesting they may have employed different rating criteria. Annotators 1 and 2 demonstrated more similar rating patterns, showing no significant differences in Openness, Extraversion, Agreeableness, and Neuroticism dimensions.

In conclusion, despite differences in the strictness of their evaluation standards, the annotators achieved good agreement in judging the relative strength of users' personality traits, particularly in the Openness, Extraversion, and Agreeableness dimensions. The assessment of Neuroticism was relatively more challenging, which aligns with our findings in Experiment 1.

G.2 Analysis of Personality Assessment Results

G.2.1 Comparison

We calculated the MAE for each of the four annotators, as presented in Table 40. We observed that

the trends are consistent with our findings in Experiment 1. Additionally, we conducted statistical analyses on the MAE for each annotator's ratings, with results shown in Tables 41, 42, 43, and 44.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	***	***	R1–R3/4/5/6 R2–R3/4/5/6 R3–R6 R4–R6 R1–R2/3/4/5/6
C	***	***	R2–R3/4/5/6 R3–R5/6 R4–R5
E	*	<i>n.s.</i>	R2–R3/4/5/6
A	<i>n.s.</i>	<i>n.s.</i>	R1–R2
N	<i>n.s.</i>	<i>n.s.</i>	R3–R4 R4–R5
AVG	***	***	R1–R3/4/5/6 R2–R3/4/5/6 R3–R5/6

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 41: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, R1-R6 = Round 1 vs Round 6.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	***	***	R1–R2/3/4/5/6 R2–R3/4/5/6 R3–R4/5/6 R4–R6 R5–R6 R1–R2/3/4/5/6
C	***	***	R2–R5/6 R3–R5 R4–R5/6 R1–R4
E	†	<i>n.s.</i>	R2–R4 R3–R4 R4–R5
A	<i>n.s.</i>	<i>n.s.</i>	None
N	<i>n.s.</i>	<i>n.s.</i>	R1–R2/4/5/6
AVG	***	***	R1–R2/3/4/5/6 R2–R3/4/5/6 R3–R4/5/6

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 42: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, R1-R6 = Round 1 vs Round 6.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	***	*	R1–R3/4/5/6 R2–R3/4/5/6 R1–R3/5/6 R2–R3/4/5/6 R3–R6 R4–R6
C	***	**	None
E	<i>n.s.</i>	<i>n.s.</i>	R1–R2/3/4/6
A	*	<i>n.s.</i>	R3–R5
N	<i>n.s.</i>	<i>n.s.</i>	R5–R6
AVG	*	<i>n.s.</i>	R2–R5/6

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 43: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, R1-R6 = Round 1 vs Round 6.

Dim.	ANOVA	Lin. Tr.	Sig. Round Pairs
O	***	***	R1–R3/4/5/6 R2–R3/4/5/6 R3–R5/6 R4–R6 R1–R2/3/4/5/6
C	***	***	R2–R3/4/5/6 R3–R5/6
E	<i>n.s.</i>	<i>n.s.</i>	R2–R3/4
A	<i>n.s.</i>	<i>n.s.</i>	R3–R4 R1–R5/6
N	**	†	R2–R5 R3–R5/6 R4–R5/6
AVG	***	***	R1–R3/4/5/6 R2–R3/4/5/6 R3–R5/6

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, † $p < 0.10$, *n.s.* = not significant

Table 44: Summary of Statistical Analysis Results for Each Dimension. Abbreviations: Dim. = Dimension, ANOVA = ANOVA Test, Lin.Tr. = Linear Trend, R1-R6 = Round 1 vs Round 6.