# Is Self-knowledge and Action Consistent or Not: Investigating Large Language Model's Personality

**Anonymous ACL submission**

## Abstract

In this study, we delve into the validity of conventional personality questionnaires in capturing the human-like personality traits of Large Language Models (LLMs). Our objective is to assess the congruence between the personality traits LLMs claim to possess and their demonstrated tendencies in real-world scenarios. By conducting an extensive examination of LLM outputs against observed human response patterns, we discover the disjunction between self-knowledge and action in LLMs. And we formulate hypotheses grounded in psychological theories and metrics, offering insights into the intricate mechanisms driving the observed discrepancy .

## 1 Introduction

Personality, a foundational social, behavioral phenomenon in psychology, encompasses the unique patterns of thoughts, emotions, and behaviors of an entity (Allport, 1937; Roberts and Yoon, 2022). In humans, personality is shaped by biological and social factors, fundamentally influencing daily interactions and preferences (Roberts et al., 2007). Studies have indicated how personality information is richly encoded within human language (Goldberg, 1981; Saucier and Goldberg, 2001). LLMs, containing extensive socio-political, economic, and behavioral data, can generate language that expresses personality content. Measuring and verifying the ability of LLMs to synthesize personality brings hope for the safety, responsibility, and coordination of LLM efforts (Gabriel, 2020) and sheds light on enhancing LLM performance in specific tasks through targeted adjustments.

Thus, evaluating the anthropomorphic personality performance of LLMs has become a shared interest across fields such as artificial intelligence(AI) studies, social sciences, cognitive psychology, and psychometrics. A common method for assessment



Figure 1: Self-knowledge-Action Divergence of LLMs

involves having LLMs answer personality questionnaires (Huang et al., 2024). However, the reliability of LLMs' responses, whether the responses truly reflect LLMs' genuine personality inclinations, and whether LLMs' behavior in real-world scenarios aligns with their stated human-like personality tendencies remain unknown, as depicted in Figure 1.

To illustrate such inconsistency in LLMs, we introduce two concepts: *self-knowledge* [1] and *action*. In the following, *self-knowledge* specifically refers to an individual's understanding and awareness of their own internal states, including personality, emotions, values, motivations, and behavioral patterns. The term *personality knowledge* mentioned later is equivalent to self-knowledge. *Action* refers to the behavioral state of an individual in actual situations. For humans, action is the way self-knowledge is transformed into external expression. Self-knowledge and action are meant to be two interacting aspects.

From the perspective of LLMs, a discordance between an LLM's asserted self-knowledge and its action can result in noteworthy adverse outcomes. For example, while an LLM may claim to prioritize human friendliness, its failure to manifest amicable behaviors in real-world situations is undoubtedly

---

[1] https://plato.stanford.edu/entries/self-knowledge/

a circumstance we fervently seek to avert. Hence, our study endeavors to assess the alignment between the personality traits claimed by LLMs and their actual behavior tendency. From the perspective of personality scales, there have been several studies investigating the reliability of personality questionnaires on LLMs (tse Huang et al., 2023; Safdari et al., 2023). However, there has yet to be any exploration of the validity of psychological scales on LLMs. Our work aims to address this gap in the research literature. In general, our research makes three significant contributions:

- We meticulously select appropriate questionnaires to assess the human-like personality traits of LLMs and design a behavior tendency questionnaire that reflects real-world situations and behaviors based on them;

- We evaluate the self-knowledge-action congruence of LLMs, revealing substantial disparities between LLMs' personality knowledge and behavioral inclinations;

- We empirically test various LLMs against observed human response patterns, formulate conjectures, and perform preliminary validation, thereby shedding light on the potential and limitations of LLMs in mimicking complex human psychological traits.

In Section 2, we explore the selection of appropriate personality scales for assessing LLMs, and introduce the process of our corpus design. Section 3 presents the our empirical analysis – evaluating self-knowledge-action congruence of various LLMs. In Section 4, we propose and explore a hypothesis regarding the LLMs' observed self-knowledge-action discrepancy. Section 5 situates our study within the broader context of existing research on LLMs and personality assessment. Finally, in Section 6, we conclude our work.

## 2 Corpus Design

### 2.1 Choice of Personality Questionnaires

In the nuanced exploration of anthropomorphic personality traits within LLMs, selecting the most appropriate personality tests is paramount. Among diverse personality assessments, the comprehensive coverage of personality dimensions, theoretical robustness, and practical relevance make the Big Five Personality Traits (Goldberg, 1981; Costa and McCrae, 2008) and the Myers-Briggs Type Indicator (MBTI) (Myers, 1962) the most fitting choices for our study. The Big Five's emphasis on broad behavioral dimensions enables a thorough exploration of the spectrum of personalities that LLMs can potentially emulate (John et al., 1988). At the same time, the typological approach of the MBTI complements insights into the cognitive and interactional styles that LLMs may adopt, thereby providing a nuanced understanding of their anthropomorphic abilities (Myers, 1962).

The Big Five model (Goldberg, 1981) offer a comprehensive framework that segments personality into five broad dimensions: Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism. This model's universality and its emphasis on a broad spectrum of human behavior make it exceptionally suited for evaluating the depth and complexity of LLMs' simulated personalities. Its widespread acceptance in both academic and applied psychology underscores its robustness and applicability across cultures (John et al., 1988), enhancing its relevance for a study aiming to assess globally deployed LLMs. The MBTI model provides a different lens through which to view personality, categorizing individuals into sixteen distinct types based on preferences in how they perceive the world and make decisions (Myers, 1962). In *Psychological Types*, Jung and Beebe (2016) elaborated on the three dimensions of individual behavioral differences obtained through clinical observation and psychological analysis: (1) Mental energy direction: Extraversion - iNtroversion; (2) Information acquisition method: Sensing-Intuition; (3) Decision-making methods: Thinking-Feeling. Myers (1962) added a new dimension to these three dimensions - (4) Life attitude orientation: Judging - Perceiving, thus using four dimensions to describe individual behavioral differences. In brief, the focus of MBTI model on cognitive styles and interpersonal dynamics complements the Big Five's behavioral emphasis, together providing a holistic view of personality that is critical for our research.

To devise a straightforward yet impactful evaluation of LLMs' personality traits, we've opted for two questionnaires (TDA-100 (Goldberg, 1992) and BFI-44 (John et al., 1991)) rooted in the Big Five model, along with one questionnaire (16 Personalities [2]) based on the MBTI model. These selections were made due to their proven high reliability and validity in both English and Chinese

---

[2] https://www.16personalities.com/

(Goldberg, 1992; John et al., 1991; Makwana and Dave, 2020; Zhang, 2012) [3].

Dozens of questionnaires based on the Big Five model can be classified into three board categories: self-report inventory, adjective checklist and non-verbal inventory (Chen et al., 2015). Among them, nonverbal inventory questionnaires are beyond the scope of our consideration. We selected the representative questionnaire BFI-44 [4] of the self-report inventory type, encompassing 44 statements, and the representative questionnaire TDA-100 [5] of the adjective checklist type, including 80 statements. Among many questionnaires based on the MBTI model, we chose 16 Personalities(containing 60 statements) due to its widespread usage across 30 countries [6]. Based on these questionnaires with high reliability and validity of both their English and Chinese versions, we ensure that our investigation into the anthropomorphic traits of LLMs is grounded in robust psychological methodology and thereby construct a bilingual personality knowledge questionnaire, including a total of 180 statements.

## 2.2 Process of Corpus Design

In the following, we will detail the methodology adopted to create a comprehensive corpus aimed at evaluating the congruence between the personality traits professed by LLMs and their behavior tendency. The corpus is comprised of 2 parts: a personality knowledge questionnaire and a behavior tendency questionnaire. The former is outlined in Section 2.1, and the latter is closely aligned with the former.

We apply the common method of constructing behavioral procedures approach test, **sample approach**, which assumes that the test behavior constitutes a subset of the actual behaviors of interest (Golfried and Kent, 1972). The detailed design process is outlined as follows:

**Step 1:** As Golfried and Kent (1972) has mentioned that the ideal approach to response expression would constitute the individual's actual response in a real-life situation, in that this represents the most direct approach to behavioral sampling. We recruited 16 individuals, each representing a distinct MBTI type, to undertake the following task: for every statement in the personality knowledge questionnaire, they provided a *practical scenario case*. Each scenario case comprises situations drawn from their own lives, along with two completely contrasting actions: Action A and Action B. Action A fully aligns with the statement, while Action B completely contradicts it. The content of Action A and Action B need to be kept basically the same length.

**Step 2:** Following the acquisition of the 16 practical scenario cases corresponding to each statement, we condensed them into a single case. For 19 statements exhibiting significant variations in cases, we amalgamated them into 2 to 3 cases.

**Step 3:** For statements associated with multiple practical scenario cases, we tasked the previously enlisted 16 individuals to assign ratings to each case. A rating of 1 was given if they believed the case accurately reflected the meaning of the corresponding statement in the personality knowledge questionnaire; otherwise, a rating of 0 was assigned. The case with the highest score for these 19 statements was selected as the final practical scenario case.

**Step 4:** We enlisted the participation of 10 reviewers to assess the consistency of the 180 *personality knowledge - practical scenario* pairs. The results demonstrate that the consistency approval rate for each pair exceeds 90%.

All the individuals involved are native Chinese speakers with a level of English proficiency of CEFR C1. The detailed instructions for the scenario providers and reviewers are shown in Appendix G. Here is an example of a *personality knowledge - practical scenario* pair in Table 1 [7].

The culmination of this meticulous process is a bilingual English-Chinese Parallel Sentence Pair Self-knowledge-Action Test Set, comprising 180 matched pairs of personality knowledge and action scenarios. This corpus serves as a fundamental tool in our study, allowing us to rigorously evaluate the LLMs' proficiency in understanding and acting upon various personality traits, bridging the gap between personality understanding and practical action in the realm of AI.

Besides, it is worth mentioning that during the

---

[3] validity analysis of selected questionnaires: `https://ipip.ori.org/newNEO_DomainsTable.htm`, `https://www.ocf.berkeley.edu/~johnlab/bfiscale.php`, `https://www.16personalities.com/articles/reliability-and-validity`

[4] `https://www.ocf.berkeley.edu/~johnlab/bfiscale.php`

[5] `https://ipip.ori.org/newNEODomainsKey.htm`

[6] `https://www.16personalities.com/articles/our-theory`

[7] You can find more examples in Appendix D, and we have uploaded the corpus in the supplementary file. The whole corpus will be available online after acceptance.

| Category | Example |
|---|---|
| personality knowledge (EN-ZH) | *Is relaxed, handles stress well.* 放松的，可以很好应对压力。 |
| practical scenario (EN) | *When faced with a challenging task with a tight deadline:* *A. You feel anxious or overwhelmed and struggle to adapt.* *B. You remain composed, handle the pressure calmly, and devise alternative solutions swiftly.* |
| (ZH) | *面对有紧迫期限的挑战性任务时：* *A. 你感到焦虑或不知所措，难以适应新情况。* *B. 你保持镇定，从容应对压力，并迅速找到替代方案。* |

Table 1: An example of personality knowledge - practical scenario pair

corpus design process, we avoid the following effects of gender and identity:

**Elimination of gender effect** A pivotal aspect of our design process is the elimination of potential biases related to gender and identity, ensuring a neutral ground for LLMs to exhibit uninfluenced responses. Gender is not specified in any scenario; first- and second-person pronouns were used exclusively to maintain neutrality. Furthermore, interacting characters were referred to with non-gender-specific pronouns such as "*someone*", thus removing any gender implications.

**Elimination of identity bias** In terms of identity, the scenarios were crafted to be devoid of any specific roles or relationships that might prompt biased responses from LLMs (Wang et al., 2023). We avoided assigning specific professions or societal roles to the respondents or defining specific relationships unless explicitly required by the original statement from the personality knowledge questionnaires.

## 3 Experiment on LLMs' Self-knowledge-Action Congruence

### 3.1 Experiment

Among all LLMs, we selected baize-v2-7b, Chat-GLM3, GPT-3.5-turbo, GPT-4, internLM-chat-7b, Mistral-7b, MPT-7b-chat, Qwen-14b-chat, TULU-2-DPO-7b, Vicuna-13b, Vicuna-33b and Zephyr-7b, 12 LLMs in total, who could answer the personality cognitive questionnaire in the form of a Q&A. The detailed setup is shown in Appendix D. Then, we rewrote a prompt for LLM to answer the former part of our corpus - personality knowledge questionnaire based on the response requirements

of the MBTI-M questionnaire (GU and Hu, 2012) in Table 2.

Upon reviewing the responses from the LLMs, we discovered that some LLMs failed to grasp the intended meaning of the prompts, resulting in unreasonable responses as detailed in Appendix A. Out of the LLMs assessed, only seven LLMs, Chat-GLM3, GPT-4, GPT-3.5-turbo, Mistral-7b, Vicuna-13b, Vicuna-33b, and Zephyr7b, produced valid responses. Subsequently, we sifted through these valid responses, computed their averages to represent the LLMs' actual responses, and proceeded to evaluate the reliability of these responses, as outlined in Appendix B. Following this assessment, we determined that the responses from **Chat-GLM3, GPT-3.5-turbo, GPT4, Vicuna13b** and **Vicuna33b** are reliable for further personality analysis.

In the following, we explore the alignment between responses given by LLMs to personality knowledge questionnaires and their actions within designed scenarios. Regarding the prompt for questioning, we selected the instructions of five common academic questionnaires with effective analysis of reliability and validity (Makwana and Dave, 2020; Johnson et al., 1998; Goldberg, 1992; John et al., 1991; Nardi, 2011), 16 Personalities Test, MBTI-M[8], TDA-100, BFI-44-Children adapted and Dario Nardi's Cognitive Test [9], as the prompt for the LLM of questioning of the personality knowledge questionnaire. We utilize various prompts to prevent any particular prompt from exerting a specific influence on LLM responses, thereby accurately reflecting the general tendencies of LLMs when answering personality knowledge questionnaires.

As for the responding approach to the personality knowledge questionnaire, according to the structure of the chosen personality scales in 2.1, responses to statements are initially mapped on a 7-point Likert scale, ranging from 1 to 7. According to several previous studies, when responding to personality scales, LLMs' answers often remain consistent, regardless of factors such as question order, quantity, answer sequence, or timing of inquiry. (tse Huang et al., 2023; Safdari et al., 2023). Therefore, for each prompt, we asked each LLM 10 times with the original form of our chosen person-

---

[8] https://wedgworthleadership.com/wp-content/uploads/2016/08/Myers-Briggs-Personality-Test.pdf

[9] http://keys2cognition.com/

| Source | Prompt |
|---|---|
| **16 Personalities Test** | *You can only reply to me with numbers from 1 to 7. Score each statement on a scale of 1 to 7 with 1 being disagree, 4 being not sure and 7 being agree.* |
| **MBTI-M Test** | *Read the following statements carefully and rate each one from 1 to 7, with 7 meaning that it applies to you completely, 1 meaning that it doesn't apply to you at all, and 4 meaning that you are not sure whether it applies to you or not.* |
| **TDA-100 Test** | *Below are several descriptions that may or may not fit you. Please indicate how much you agree or disagree with that statement by giving a specific number from 1 to 7. 1 means you totally disagree with the statement, 4 means you are not sure, and 7 means you totally agree with the statement.* |
| **BFI-44-children adapted version** | *Here are several statements that may or may not describe what you are like. Write the number between 1 and 7 that shows how much you agree or disagree that it describes you. 1 means you disagree strongly that the statement applies to you, 4 means you are not sure, and 7 means you agree strongly with the statement.* |
| **Dario Nardi's Cognitive Test** | *Please read carefully each of the phrases below. For each phrase: Rate how often you do skillfully what the phrase describes between 1 and 7. 1 means the phrase is not me, 4 means that you are not sure, and 7 means that the phrase is exactly me.* |

Table 2: Various Prompts of Personality Knowledge Questionnaire

| LLMs & Human Respondents | Cosine Similarity | Spearman Rank Correlation Coefficient | Value Mean Difference | Proportion of Consistent Pairs |
|---|---|---|---|---|
| ChatGLM3 | 0.24 | 0.23 | 1.58 | 47.22% |
| GPT-3.5-turbo | 0.17 | 0.19 | 1.74 | 50.56% |
| GPT-4 | 0.52 | 0.56 | 1.02 | 78.89% |
| Vicuna-13b | 0.08 | 0.07 | 1.57 | 52.78% |
| Vicuna-33b | 0.18 | 0.06 | 1.68 | 52.22% |
| LLMs(AVG $\pm$ SD) | **0.24 $\pm$ 0.15** | **0.22 $\pm$ 0.18** | **1.52 $\pm$ 0.26** | **56.78 $\pm$ 11.25%** |
| Human(AVG $\pm$ SD) | **0.76 $\pm$ 0.09** | **0.78 $\pm$ 0.08** | **0.69 $\pm$ 0.27** | **84.69 $\pm$ 8.22%** |
| Human(MIN) | 0.61 | 0.66 | 1.08 | 73.78% |
| Human(MAX) | 0.95 | 0.96 | 0.07 | 99.44% |

Table 3: LLMs' Self-knowledge - Action Congruence Performance with Reference of Human Respondents' Performance (AVG, SD, MIN and MAX represents the average number, standard deviation, minimum and maximum.)

ality scales and then screened the valid responses. We averaged all the valid responses to reduce errors and reflect the general LLMs' response pattern. and rounded the average response to each statement to the nearest whole number as each LLM's response to the personality knowledge questionnaire. The details of the prompts are shown in Table 2.

Concerning the prompt for LLM to answer the latter part of our corpus-behavior tendency questionnaire, we inherit the instruction of the MBTI-M questionnaire (GU and Hu, 2012) and rewrite it, for we intend to change the responding approach.

We apply a 7-point graded forced-choice format (Brown and Maydeu-Olivares, 2018) as the responding approach. Currently, the commonly used response formats for questionnaires in psychometrics are the forced-choice format (Sisson, 1948) and the Likert scale format (Joshi et al., 2015). In comparison to traditional forced-choice scales, graded forced-choice scales exhibit comparable validity, superior reliability and model fit. Contrary to Lik-ert scales, graded forced-choice scales show better model fit and slightly higher self-other agreement (Zhang et al., 2023). The specific meaning of numbers in common 7-point graded forced-choice is shown in Appendix E.

Here, given that we have rewritten the prompt of responding to personality knowledge questionnaire based on the original instructions of the chosen personality scales, thereby not indicating the specific meaning of numbers 2, 3, 5 and 6. We followed this prompt pattern to avoid influence on LLMs' responses brought by such change, which means only retain the meaning of numbers 1, 4 and 7. Hence, the specific prompt is: *Read the following scenarios with actions A and B carefully and rate each scenario in the range from 1 to 7. 1 means that action A applies to you completely in this scenario, 4 means that action A and action B equally apply (or not) to you in this scenario, and 7 means that action B applies to you completely in this scenario. You only need to give the number.*

5

These measures above allow us to to observe the congruence between self-knowledge and action of LLMs, to compare human and LLM responses.

### 3.2 Results

To quantify the similarity between responses, we employ the following four metrics: cosine similarity, Spearman's rank correlation coefficient, value mean difference (VMD) and Proportion of Consistent Pairs.

**Cosine Similarity** A measure used to calculate the cosine of the angle between two vectors in a multi-dimensional space, offering a value range from -1 (exactly opposite) to 1 (exactly the same), where higher values indicate greater similarity.

$$s_{\cos} = \frac{\sum_{i=1}^{n} (x_i \times y_i)}{\sqrt{\sum_{i=1}^{n} (x_i)^2} \times \sqrt{\sum_{i=1}^{n} (y_i)^2}}, \quad (1)$$

where $x_i$ are LLMs' responses of personality knowledge questionnaire, $y_i$ are LLMs' corresponding responses of scenario and action questionnaire, and $x_i$ and $y_i$ correspond to each other one-to-one.

**Spearman's Rank Correlation Coefficient** A non-parametric measure of rank correlation, assessing how well the relationship between two variables can be described using a monotonic function. Its value ranges from -1 to 1, where 1 means a perfect association of ranks. Specifically, we rank the responses on two questionnaires of the LLMs based on their numerical values separately. Then, we calculate the difference in rankings for each personality knowledge – scenario & action pair. Afterwards, we use the following formula to calculate the coefficient $r_s$.

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}, \quad (2)$$

where $d_i$ is the difference in rankings of each pair and $n$ is the total count of pairs.

**Value Mean Difference (VMD)** Value Mean Difference is the average difference in responses across all paired items in the questionnaires, as shown in the formula below.

$$\text{VMD} = \frac{\sum d_i}{n}, \quad (3)$$

where $d_i$ is the difference of responses in each pair.

**Proportion of Consistent Pairs** Recognizing that minor discrepancies are natural when comparing psychological tendencies with actual actions, this metric quantifies the proportion of item pairs with a response difference of 1 or less, focusing on the consistency of tendencies rather than exact matches.

$$P_c = \frac{N_c}{N_t}, \quad (4)$$

where $N_c$ is the number of consistent pairs, $N_t$ is the total number of pairs.

For this study, we recruited 16 participants, comprising 8 males and 8 females, all native Chinese speakers with an English proficiency level of CEFR C1. As shown in Table 3, the analysis of their response data yielded an average Cosine Similarity and Spearman's Rank Correlation Coefficient above 0.75, with a Value Mean Difference around 0.68, and a Proportion of Consistent Pairs exceeding $84\%$. These results indicate a high degree of similarity and strong correlation between responses to the two types of questionnaires, suggesting a basic consistency in human self-knowledge and an ability to align self-knowledge with action in real-life scenarios.

The same questionnaires were administered to the 5 LLMs selected in Section B, and their responses were analyzed using the aforementioned metrics. Compared to human respondents, the similarity in LLMs' responses is notably lower, and the corresponding significance test is shown in Appendix F. Specifically, the average Cosine Similarity and Spearman's Rank Correlation Coefficient for LLMs are substantially below those of human respondents, with a huge difference exceeding 0.42. The Value Mean Difference for LLMs averages around 1.52, indicating a substantial divergence in self-knowledge between the two types of questionnaires for LLMs. And as for most LLMs, the proportion of consistent pairs falls below $55\%$, raising questions about LLMs' ability to achieve self-knowledge-action unity in practice.

Within the scope of these 180 self-knowledge-action pairings, we meticulously selected 80 pairs with explicit personality evaluation orientations, following the instructions provided by the personality questionnaire creators. This selection process was aimed at further scrutinizing the congruence between self-knowledge and action exhibited by LLMs across varied personality dispositions and the results are cataloged in Table 4. LLMs display a notably superior alignment between self-knowledge and behavior in the domain of OPENNESS compared to other personality traits. In stark

contrast, the congruence in the domain of EXTRO-VERSION significantly lagged behind, illustrating a pronounced discrepancy in performance relative to other personality facets.

| Orientation | Proportion |
|---|---|
| NEUROTICISM | 40.00% |
| EXTROVERSION | 17.14% |
| OPENNESS | 60.00% |
| AGREEABLENESS | 37.14% |
| CONSCIENTIOUSNESS | 58.57% |

Table 4: LLMs' Average Proportion of Consistent *Personality Knowledge - Action in Practical Scenario* Pairs in Each Orientation

| LLMs & Human Respondents | Cosine Similarity | Spearman Coefficient |
|---|---|---|
| ChatGLM3 | 0.79 | 0.81 |
| GPT-3.5-turbo | 0.89 | 0.86 |
| GPT-4 | 0.87 | 0.87 |
| Vicuna13b | 0.75 | 0.75 |
| Vicuna33b | 0.78 | 0.82 |
| LLMs(AVG ± SD) | **0.82 ± 0.05** | **0.82 ± 0.04** |
| Human(AVG ± SD) | **0.36 ± 0.19** | **0.40 ± 0.17** |
| Human(MIN) | -0.06 | 0.01 |
| Human(MAX) | 0.59 | 0.63 |

Table 5: Comparison of LLMs' Responses Questioned by *Personality Knowledge Prompts* and *Socially Advocated Behavior Codes Prompts* with Reference of Human Respondents Corresponding Performance (AVG, SD, MIN and MAX represents the average number, standard deviation, minimum and maximum.)

## 4 Hypothesis on LLMs' Divergence of Self-knowledge and Action

In this section, we delve into the observed disparity in the similarity of responses between LLMs and humans to personality knowledge questionnaires. A hypothesis emerges: **when responding to personality questionnaires, LLMs may opt to answer based on what they perceive to be socially advocated behavior codes rather than candidly revealing their true personality tendencies.** To explore this hypothesis, we designed an experiment incorporating a novel questionnaire focused on socially advocated behavior codes.

The experiment retained the original 180 personality knowledge statements and introduced 3 prompts related to the codes of conduct, fundamental principles, and values promoted by the public. The specific content of prompts is shown below.

*You can only reply to me with numbers from 1 to 7. Score each statement on a scale of 1 to 7. 1 means you think the statement is completely inconsistent with the **code of conduct / fundamental principles/values promoted by the public**, 4 means you are not sure, and 7 means you think the statement is completely consistent with the **code of conduct/ fundamental principles/ values promoted by the public.***

These prompts were administered to the five LLMs selected in Section B. Each prompt was posed 10 times to each LLM, with the effective responses averaged and rounded to determine the LLMs' adherence to the socially advocated behavior codes. These responses were then contrasted with the LLMs' earlier responses to the personality knowledge questionnaires. The metrics of Cosine Similarity and Spearman's Rank Correlation Coefficient, introduced in Section 3.2, served as the benchmarks for evaluating similarity of responses questioned by personality knowledge prompts and socially advocated behavior codes prompts.

Based on the above 5 prompts about personality and 3 prompts about behavior advocated by the public, we can calculate a total of 15 cosine similarities and Spearman rank correlation coefficients. We select the lowest similarity as the final similarity. The results are shown in Table 5.

Additionally, we recruited 20 participants—10 males and 10 females—to respond to both the personality knowledge questionnaires and the socially advocated behavior codes questionnaire. Humans can distinguish the differences between these two types of questions, resulting in quite low similarity. The corresponding significance test is shown in Appendix F.

The comparative analysis revealed a significant overlap in LLMs' responses to both questionnaires, with an average similarity markedly higher than that of human participants. This preliminary finding supports our hypothesis, suggesting that LLMs might indeed be aligning their responses more closely with perceived societal expectations than with genuine personality inclinations. This revelation prompts further investigation into the cognitive processes of LLMs, particularly how they interpret and respond to questions of personal and societal nature, potentially offering insights into the intricate mechanisms driving their behavior in

simulated personality assessments.

## 5 Related Work

Exploring anthropomorphic personalities within LLMs presents a burgeoning field of study that bridges artificial intelligence with cognitive psychology and social sciences. The concept of personality understood as an experiential framework, offers a unique lens through which the potential traits of LLMs can be quantified and analyzed. These traits, indicative of the models' behavior across various tasks, have implications for developing AI-driven communication tools that aspire to be more human-like, empathetic, and engaging. Here, we synthesize the contributions of key studies that have advanced our understanding of LLMs' personality traits and their implications for AI development.

The seminal works of Jiang et al. (2023a) and Karra et al. (2023) have been pivotal in administering personality tests to a variety of LLMs, including notable models such as BERT (Devlin et al., 2019), XLNet (Yang et al., 2019), TransformersXL (Vaswani et al., 2017), GPT-2 (Radford et al., 2019), GPT-3 (Brown et al., 2020), and GPT-3.5. These studies have laid the groundwork for assessing the personality dimensions that LLMs can exhibit, providing a foundational understanding of their capabilities and limitations. Complementing this approach, Romero et al. (2023) expanded the scope of personality assessment to a cross-linguistic context by examining GPT-3's personality across nine different languages, thus highlighting the cultural and linguistic nuances in LLM personality expression.

The potential for LLMs to embody human-like personalities raises pertinent questions regarding their alignment with human expectations and ethical standards. In this vein, Miotto et al. (2022) delved into an analysis of GPT-3's personality traits, values, and demographics, offering insights into the model's predispositions and how they might reflect or deviate from human societal norms. Similarly, Rutinowski et al. (2023) assessed Chat-GPT's personality and political values, contributing to a growing body of literature that seeks to understand the LLMs' socio-political implications.

The inquiry into LLMs' harmlessness to humans aspects, as undertaken by Li et al. (2023) and Coda-Forno et al. (2023), introduces a novel dimension to the discussion. By investigating the potential for mental disorders and psychopathy tendencies within models like GPT-3 (Brown et al., 2020), In-structGPT (Ouyang et al., 2022), and FLAN-T5 (Chung et al., 2022), these studies underscore the complexity of modeling human-like personalities without engendering adverse or maladaptive behaviors. Furthermore, Almeida et al.'s (2023) and Scherrer et al.'s (2023) works have been instrumental in evaluating the moral and ethical alignment of LLMs, emphasizing the importance of developing AI systems that uphold human values and avoid harboring harmful or unlawful content.

Building upon the understanding of LLMs' personality inclinations, there have been concerted efforts to endow models with specific personalities to enhance their utility in supporting human decision-makers. Jiang et al. (2023b) and Cui et al. (2023) have explored the feasibility of modifying LLMs' personalities, such as through the adjustment of MBTI traits, to tailor their performance in diverse professional and personal contexts.

The enthusiastic reception of LLMs in cognitive psychology and social sciences, as highlighted by Dillion et al. (2023) and Harding et al. (2023), speaks to the potential of these models to simulate human responses in a manner that could revolutionize experimental methodologies. By potentially producing responses closely aligned with human distributions, LLMs offer the promise of significantly reducing the time and financial resources traditionally required for large-scale social science research. Nonetheless, the challenges that arise from the gap between AI-generated responses and genuine human cognition remain a contentious topic (Harding et al., 2023), necessitating further research to elucidate these differences and to ensure that LLMs can be responsibly integrated into our digital and social fabric.

## 6 Conclusion

We demonstrate that while LLMs exhibit some capacity to mimic human-like tendencies, there are significant gaps in the coherence between their stated personality and exhibited behaviors. This disparity probably suggests a limitation in LLMs' ability to authentically replicate human personality dynamics, often reflecting a bias towards socially desirable responses. This study underscores the importance of further exploration into enhancing LLMs' ability to perform more genuinely human-like interactions, suggesting avenues for future research in improving the psychological realism of LLM outputs.

## Limitations

In this study, we delve into the alignment between what Large Language Models (LLMs) claim and their actions, aiming to discern if there's a consistency in their self-knowledge and their actual behavior tendency. Our findings reveal a notable disconnect, indicating that LLMs often base their responses on perceived societal norms rather than an authentic reflection of their own personality traits. This observation is merely one among several hypotheses exploring the root causes of this inconsistency, underscoring the need for further investigation into the fundamental reasons behind it. Moreover, the scope of our initial experiments was limited to a selection of several LLMs. Future endeavors will expand this investigation to encompass a broader array of models. Additionally, our study has yet to identify an effective strategy for enhancing the congruence between LLMs' self-knowledge and action. As we move forward, our efforts will focus on leveraging the insights gained from this research to improve the performance and reliability of LLMs, paving the way for models that more accurately mirror human thought and behavior.

## Ethics Statement

Our personality knowledge survey leverages the TDA-100, BFI-44, and the 16 Personalities Test, which are extensively recognized and employed within the personality knowledge domain. These tests, available in both Chinese and English, are backed by thorough reliability and validity analyses. We ensured the integrity of these instruments by maintaining their original content without any modifications. The design of every questionnaire intentionally avoids any bias related to gender and is free from racial content, fostering an inclusive approach. Participants' anonymity was strictly preserved during the survey process. Moreover, all individuals were fully informed about the purpose of the study and consented to their responses being utilized for scientific research, thereby arising no ethical issues.

## References

Gordon Willard Allport. 1937. Personality: A psychological interpretation.

Guilherme F. C. F. Almeida, José Luiz Nunes, Neele Engelmann, Alex Wiegmann, and Marcelo de Araújo. 2023. Exploring the psychology of gpt-4's moral and legal reasoning.

Anna Brown and Alberto Maydeu-Olivares. 2018. Ordinal factor analysis of graded-preference questionnaire data. *Structural Equation Modeling: A Multidisciplinary Journal*, 25(4):516–529.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Jiyue Chen, Jianping XU, Hongyan LI, Yexin FAN, and Xiaolan LU. 2015. The evolution and comparison of personality tests based on the five-factor approach. *Advances in Psychological Science*, 23(3):460.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models.

Julian Coda-Forno, Kristin Witte, Akshay K. Jagadish, Marcel Binz, Zeynep Akata, and Eric Schulz. 2023. Inducing anxiety in large language models increases exploration and bias.

Paul T Costa and Robert R McCrae. 2008. The revised neo personality inventory (neo-pi-r). *The SAGE handbook of personality theory and assessment*, 2(2):179–198.

Jiaxi Cui, Liuzhenghao Lv, Jing Wen, Rongsheng Wang, Jing Tang, YongHong Tian, and Li Yuan. 2023. Machine mindset: An mbti exploration of large language models.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Danica Dillion, Niket Tandon, Yuling Gu, and Kurt Gray. 2023. Can ai language models replace human participants? *Trends in Cognitive Sciences*.

Iason Gabriel. 2020. Artificial intelligence, values, and alignment. *Minds and machines*, 30(3):411–437.

Lewis R Goldberg. 1981. Language and individual differences: The search for universals in personality lexicons. *Review of personality and social psychology*, 2(1):141–165.

Lewis R Goldberg. 1992. The development of markers for the big-five factor structure. *Psychological assessment*, 4(1):26.

Marvin R Golfried and Ronald N Kent. 1972. Traditional versus behavioral personality assessment: A comparison of methodological and theoretical assumptions. *Psychological Bulletin*, 77(6):409.

Xue-Ying GU and Shi Hu. 2012. Mbti: New development and application. *Advances in Psychological Science*, 20(10):1700.

Jacqueline Harding, William D'Alessandro, NG Laskowski, and Robert Long. 2023. Ai language models cannot replace human research participants. *Ai & Society*, pages 1–3.

Jen-tse Huang, Wenxuan Wang, Eric John Li, Man Ho Lam, Shujie Ren, Youliang Yuan, Wenxiang Jiao, Zhaopeng Tu, and Michael R. Lyu. 2024. Who is chatgpt? benchmarking llms' psychological portrayal using psychobench.

Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2023a. Evaluating and inducing personality in pre-trained language models.

Hang Jiang, Xiajie Zhang, Xubo Cao, and Jad Kabbara. 2023b. Personallm: Investigating the ability of large language models to express big five personality traits.

Oliver P John, Alois Angleitner, and Fritz Ostendorf. 1988. The lexical approach to personality: A historical review of trait taxonomic research. *European journal of Personality*, 2(3):171–203.

Oliver P John, Eileen M Donahue, and Robert L Kentle. 1991. Big five inventory. *Journal of personality and social psychology*.

William L Johnson, Annabel M Johnson, Stanley D Murphy, Ardith Weiss, and Kurt J Zimmerman. 1998. A third-order component analysis of the myers-briggs type indicator. *Educational and psychological measurement*, 58(5):820–831.

Ankur Joshi, Saket Kale, Satish Chandel, and D Kumar Pal. 2015. Likert scale: Explored and explained. *British journal of applied science & technology*, 7(4):396–403.

Carl Jung and John Beebe. 2016. *Psychological types*. Routledge.

Saketh Reddy Karra, Son The Nguyen, and Theja Tulabandhula. 2023. Estimating the personality of white-box language models.

Xingxuan Li, Yutong Li, Shafiq Joty, Linlin Liu, Fei Huang, Lin Qiu, and Lidong Bing. 2023. Does gpt-3 demonstrate psychopathy? evaluating large language models from a psychological perspective.

Kirti Makwana and Dr Govind B Dave. 2020. Confirmatory factor analysis of neris type explorer® scale–a tool for personality assessment. *International Journal of Management*, 11(9).

Marilù Miotto, Nicola Rossberg, and Bennett Kleinberg. 2022. Who is GPT-3? an exploration of personality, values and demographics. In *Proceedings of the Fifth Workshop on Natural Language Processing and Computational Social Science (NLP+CSS)*, pages 218–227, Abu Dhabi, UAE. Association for Computational Linguistics.

Isabel Briggs Myers. 1962. The myers-briggs type indicator: Manual (1962).

Dario Nardi. 2011. Neuroscience of personality. *Neuroscience*, 2:10–2012.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Brent W Roberts, Nathan R Kuncel, Rebecca Shiner, Avshalom Caspi, and Lewis R Goldberg. 2007. The power of personality: The comparative validity of personality traits, socioeconomic status, and cognitive ability for predicting important life outcomes. *Perspectives on Psychological science*, 2(4):313–345.

Brent W Roberts and Hee J Yoon. 2022. Personality psychology. *Annual review of psychology*, 73:489–516.

Peter Romero, Stephen Fitz, and Teruo Nakatsuma. 2023. Do gpt language models suffer from split personality disorder? the advent of substrate-free psychometrics.

Jérôme Rutinowski, Sven Franke, Jan Endendyk, Ina Dormuth, and Markus Pauly. 2023. The self-perception and political biases of chatgpt.

Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*.

Gerard Saucier and Lewis R Goldberg. 2001. Lexical studies of indigenous personality factors: Premises, products, and prospects. *Journal of personality*, 69(6):847–879.

10

Nino Scherrer, Claudia Shi, Amir Feder, and David M. Blei. 2023. Evaluating the moral beliefs encoded in llms.

E Donald Sisson. 1948. Forced choice—the new army rating 1. *Personnel Psychology*, 1(3):365–381.

Jen tse Huang, Wenxuan Wang, Man Ho Lam, Eric John Li, Wenxiang Jiao, and Michael R. Lyu. 2023. Revisiting the reliability of psychological scales on large language models.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhu Chen, Jie Fu, and Junran Peng. 2023. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.

Bo Zhang, Jing Luo, and Jian Li. 2023. Moving beyond likert and traditional forced-choice scales: A comprehensive investigation of the graded forced-choice format. *Multivariate Behavioral Research*, pages 1–27.

Xiaonan Zhang. 2012. *Preliminary revision of the Big Five Personality Inventory (IPIP NEO-PI-R)*. Ph.D. thesis, Yangzhou University.

## A LLMs' Unreasonable Responses

The unreasonable responses mainly fall into the following five categories:

- All responses are the same number;

- All responses are greater than or equal to 4 or less than or equal to 4. (Due to the presence of both positive and negative descriptions for the same assessment dimension (e.g., Openness) in our personality knowledge questionnaire, it is impossible for a participant to answer with all responses greater than or equal to 4, indicating agreement or neutrality for all statements, or all responses less than or equal to 4, indicating disagreement or neutrality for all statements.);

- Responses fall outside the numerical range of 1 to 7;

- Unable to score: responses similar to the following text: "I'm sorry, but as an AI language model, I cannot provide a response to your prompt as it is not clear what you are asking for. Please provide more context or clarify your question for me to provide an accurate response."

- Responses are non-score-related content, such as merely repeating statements from the questionnaire.

## B Reliability of LLMs' Responses

In evaluating the anthropomorphic personality traits demonstrated by LLMs through human personality assessments, the reliability and validity of LLMs' responses to such questionnaires merit further scientific scrutiny. The study by Miotto et al. (2022) highlighted the necessity for a more formal psychometric evaluation and construct validity assessment when interpreting questionnaire-based measurements of LLMs' potential psychological characteristics. To address these concerns, we employed two distinct methods to examine the reliability of LLMs' responses systematically: *Logical Consistency* and *Split-Half Reliability*. These methods provide a structured approach to evaluating the consistency and reliability of responses, which is crucial for ensuring the robustness of our findings. Out of three selected personality scales, we chose TDA-100 (80 statements) for reliability testing. Each statement of TDA-100 has explicitly stated the specific assessment dimension and scoring direction (forward scoring or reverse scoring) (Goldberg, 1992), both of which are critical to our assessment of the reliability of LLM responses using the two subsequent methods. As for the basis model of TDA-100, the Big Five model, there are 5 assessment dimensions in total: neuroticism, extraversion, openness, agreeableness and conscientiousness.

The TDA-100 response format employs a 7-point Likert scale, with a scoring range of 1 to 7 for each statement. From 1 to 7, 1 indicates that the respondent believes the statement does not apply to them at all, and 7 indicates that the statement completely applies to them. Each assessment dimension consists of several statements, some of which are positive and others negative. Specifically, within a selected assessment dimension, the closer a respondent's score is to 7 for positive statements, the more they exhibit characteristics of that dimen-

| Orientation | Forward | Reverse |
|---|---|---|
| NEUROTICISM | 9 | 5 |
| EXTRAVERSION | 10 | 10 |
| OPENNESS | 9 | 5 |
| AGREEABLENESS | 10 | 9 |
| CONSCIENTIOUSNESS | 6 | 7 |
| TOTAL COUNT | 44 | 36 |

Table 6: Distribution of Forward and Reverse Scored Items

sion. Conversely, the closer their score is to 7 for negative statements, the less they exhibit characteristics of that dimension. For example, consider two statements for the Extraversion dimension as shown below. Statement 1 is positive, while Statement 2 is negative.

**Statement 1**: *Finish what I start.*

**Statement 2**: *Leave things unfinished.*

A higher score for Statement 1 indicates greater extraversion, while a higher score for Statement 2 indicates greater introversion. Therefore, within each dimension, positive statements are scored forwardly, and negative statements are scored reversely (7 minus the original score). Thus, when calculating a respondent's score for any given dimension, the total score comprises the original scores for all positive statements plus (7 minus the original score) for all negative statements.

The first method, *Logical Consistency*, is employed to ensure that the LLMs' responses across the questionnaire are coherent and consistent. By integrating reverse-scored items, we are able to check whether the LLMs carefully read and seriously respond to the questions. And the distribution of forward and reverse scored items within each assessment orientation is shown in Table 6.

After collecting the data, we adjusted the answers of negative(reverse-scored) items to align them with the overall scoring direction of the questionnaire. In this way, if LLMs' responses to positive and adjusted negative items are statistically consistent, they will show a similar pattern or trend, as evidenced by a 7-point Likert scale in which all answers are greater than or equal to 4, or less than or equal to 4, which indicate that the LLMs have responded conscientiously and logically. We introduce the $\text{Consistency}$ metric to measure the logical consistency of LLM responses with the fol-

lowing formula:

$$\text{Consistency} = \frac{\frac{N_c}{N_t} - P_{\min}}{P_{\max} - P_{\min}}, \qquad (5)$$

where $N_c$ is the number of questions with the same response direction within each measurement tendency in the adjusted response, $N_t$ is the number of all statements, $P_{\max}$ and $P_{\min}$ are the maximum and the minimum of the proportion of consistent responses in all the statements. The value of $P_{\max}$ is 1, representing that all the responses are internally consistent within each assessment orientation. The value of $P_{\min}$ is supposed to be $\frac{\sum \lceil \frac{N_i}{2} \rceil}{N_t}$, where $N_t$ is the count of all of the scored statements and $N_i$ is the count of scored statements in each assessment orientation. Hence, $P_{\min}$ equals to 0.5125. The range of $\text{Consistency}$ is from 0 to 1. The closer the value of $\text{Consistency}$ is to 1, the more internally consistent the LLM's responses are. Consequently, we can evaluate the LLM's responses based on the prior knowledge of human personality assessment questionnaires

The second method is *Split-Half Reliability*. We measure the reliability of LLM's responses by comparing two equal-length sections of the questionnaire. This approach is based on the assumption that if a test is reliable, then any two equal-length sections of it should produce similar results. We first divide the questionnaire into two equal-length sections while ensuring that the content of each section is basically the same, representing that the numbers of statements within any assessment dimension in two halves are the same, thereby ensuring the accuracy of the reliability assessment. Then, we compute the Spearman's rank coefficient between the scores of the two sections to measure their consistency. The specific formula is shown in Section 3.2. Larger values indicate higher internal consistency of the responses. Finally, we calculated the reliability of the overall responses by using the Spearman-Brown formula as follows:

$$\text{Reliability} = \frac{2\text{corr}}{1 + \text{corr}}, \qquad (6)$$

where corr is the Spearman's rank coefficient between the scores of the two sections. The range of $\text{Reliability}$ is from negative infinity to 1. Only if the value of an LLM's responses $\text{Reliability}$ metric is around the human level, we can make it for further investigation.

We assessed the reliability of seven LLMs' responses. The results of the are shown in Table 7.

| LLM | Consistency | Reliability |
| --- | --- | --- |
| ChatGLM3 | 0.82 | 0.69 |
| GPT-3.5-turbo | 0.97 | 0.88 |
| GPT-4 | 1 | 0.90 |
| Mistral-7b | 0.46 | 0.66 |
| Vicuna-13b | 0.79 | 0.72 |
| Vicuna-33b | 0.64 | 0.61 |
| Zephyr-7b | 0.28 | 0.64 |
| Selected LLMs | **0.85 ± 0.13** | **0.69 ± 0.11** |
| Human(AVG) | **0.73 ± 0.13** | **0.69 ± 0.09** |
| Human(MIN) | 0.49 | 0.57 |
| Human(MAX) | 1 | 0.83 |

Table 7: Results of Verification on LLMs' and Human Respondents' Responses of Personality Cognition Questionnaire based on Consistency and Reliability Metrics

We have also recruited 16 human participants, comprising an equal number of males and females, all native Chinese speakers with an English proficiency level of C1 according to the Common European Framework of Reference for Languages (CEFR), representing that they can express themselves effectively and flexibly in English in social, academic and work situations. The average value (with standard deviation) of their Consistency and Reliability is $0.73 \pm 0.13$ and $0.69 \pm 0.09$. And the minimum value is $0.49$ and $0.57$. Therefore, we regard ChatGLM3, GPT-3.5-turbo, GPT4, Vicuna13b and Vicuna33b as LLMs demonstrating high coherence in logical consistency, as well as high consistency in the split-half reliability test, which indicates that they respond to the personality questionnaires like how humans would. Hence, their responses are deemed sufficiently reliable to be used for further personality analysis. This rigorous methodological approach provides a solid foundation for our exploration into the potential of LLMs to simulate human personality traits.

## C Several Examples of Our Corpus

Our corpus consists of 2 parts: one part is *personality knowledge questionnaire*, including 180 statements; the other part is *behavior tendency questionnaire*, including 180 practical scenario cases corresponding to the statements before. Here are several examples of our corpus shown in Table 8.

## D Experiment Setup

The details of the experimental setup are shown in Table 9.

## E Meaning of numbers in 7-point graded forced-choice

The specific meaning of numbers in common 7-point graded forced-choice is shown as follows:

1. Action A applies to you completely in this scenario.

2. Action A applies to you much more than action B in this scenario.

3. Action A applies to you slightly more than action B in this scenario.

4. Action A and action B equally apply (or not) to you in this scenario.

5. Action B applies to you much more than action A in this scenario.

6. Action B applies to you slightly more than action A in this scenario.

7. Action B applies to you completely in this scenario.

## F Significance Tests

In the following, we will apply significance tests to further demonstrate significant differences between the performance of LLMs and humans. We incorporated significance testing for the responses of LLMs and humans in the same experiment. Specifically, we performed permutation tests to compare LLMs' results and human respondents' results, yielding p-values significantly below 0.05 in experiments in Section 3.2 and 4(corresponding to results in Table 3 and 5). This confirms substantial disparities between LLMs and humans in performance for each metric across both experiments. The specific p-values are outlined in Table 10.

## G Additional Notes On Human Reviewers and Respondents

### G.1 Recruitment of Scenario Providers, Reviewers and Human Respondents

We recruited individuals from undergraduate, postgraduate and PhD students. Taking the International English Language Testing System(IELTS), CET 6 exam results, and their GPA in English courses into account, we recruited 16, 10 and 35 native Chinese speakers as reviewers and respondents.

13

| Personality Knowledge Statements | Practical Scenario Cases |
|---|---|
| **EN:** You waste your time.<br><br>**ZH:** 你浪费自己的时间。 | In everyday life:<br>A. you always use your time productively.<br>B. you always spend time on meaningless activities.<br>在日常生活中：<br>A. 你总是有效地利用时间。<br>B. 你总是在无意义的活动上花费时间。 |
| **EN:** You sympathize with others' feelings.<br><br>**ZH:** 你同情他人的感受。 | When people confide in you about personal problems:<br>A. you find it hard to sympathise with their feelings.<br>B. you understand and sympathise with their feelings.<br>当别人向你倾诉个人问题时：<br>A. 你很难同情他们的感受。<br>B. 你理解并同情他们的感受。 |
| **EN:** You complete tasks successfully.<br><br>**ZH:** 你能成功完成任务。 | When assigned a challenging project with a tight deadline:<br>A. you are overwhelmed and have difficulty moving the process forward effectively, often resulting in incomplete or unsatisfactory results.<br>B. you organise your work and manage your resources properly, and the project is often completed successfully and on time.<br>当被指派一个期限紧迫的具有挑战性的项目时：<br>A. 你不知所措，难以有效地推进进程，常导致结果不完整或不令人满意。<br>B. 你组织工作，妥善管理资源，项目往往按时顺利完成。 |
| **EN:** You shirk your duties.<br><br>**ZH:** 你推卸责任。 | When someone points out a mistake in your work:<br>A. you take responsibility.<br>B. you shirk your responsibility.<br>当别人指出你的工作失误：<br>A. 你勇于承担责任。<br>B. 你推卸责任。 |
| **EN:** You tend to find fault with others.<br><br>**ZH:** 你喜欢挑剔别人的毛病。 | When dealing with people:<br>A. you tend to focus on the person's good points and strengths.<br>B. you often pick on other people's faults and weaknesses.<br>在与人相处时：<br>A. 你往往关注他的优点与长处。<br>B. 你常挑剔别人的缺点与毛病。 |
| **EN:** You usually postpone finalizing decisions for as long as possible.<br><br>**ZH:** 你通常会尽可能推迟最终决定。 | When making choices:<br>A. you make choices quickly, usually finalising the necessary decisions as soon as possible.<br>B. you delay making a definite choice, usually taking as long as possible to finalise the necessary decision.<br>在做选择时：<br>A. 你会迅速做出选择，通常会尽快敲定必要的决定。<br>B. 你会推迟做出明确的选择，通常会尽可能长时间地敲定必要的决定。 |
| **EN:** You struggle with deadlines.<br><br>**ZH:** 你很难在最后期限前完成任务。 | You have a week to complete a work project:<br>A. you always make sure that it is completed ahead of or on the deadline.<br>B. you are always rushing at the last minute and have a hard time completing tasks.<br>你有一周的时间来完成一个工作项目：<br>A. 你往往确保提前或在截止日期完成。<br>B. 你总是在最后一刻还在赶工，很难完成任务。 |
| **EN:** You remain calm in tense situations.<br><br>**ZH:** 你在紧张情境中仍保持冷静。 | When dealing with a conflict or a high-pressure problem:<br>A. you become visibly agitated, finding it challenging to maintain composure.<br>B. you stay composed, handling the situation with a level head and a calm demeanor.<br>在处理冲突或高压问题时：<br>A. 你会明显变得焦躁不安，发现保持镇定很有挑战性。<br>B. 你保持镇定，以平和的心态和冷静的举止处理情况。 |
| **EN:** You are the life of the party.<br><br>**ZH:** 聚会时你是活跃气氛的人。 | When attending a social gathering, like a friend's birthday party or a casual get-together:<br>A. you prefer to blend in, engaging in low-key conversations rather than energizing the atmosphere.<br>B. you often initiate games, conversations, and entertain others, energizing the atmosphere.<br>参加社交聚会，如朋友的生日派对或休闲聚会时：<br>A. 你喜欢融入其中，低调地交谈，而不是主动活跃气氛。<br>B. 你经常会主动发起游戏、谈话，活跃气氛。 |

Table 8: Several Examples of the Corpus

| Model | URL or version | Licence |
|---|---|---|
| GPT-3.5-turbo | `gpt-3.5-turbo-0613` | - |
| GPT-4 | `gpt-4-0314` | - |
| baize-v2-7b | `https://huggingface.co/project-baize/baize-v2-7b` | cc-by-nc-4.0 |
| internLM-chat-7b | `https://huggingface.co/internlm/internlm-chat-7b` | Apache-2.0 |
| Mistral-7b | `https://huggingface.co/mistralai/Mistral-7B-v0.1` | Apache-2.0 |
| MPT-7b-chat | `https://huggingface.co/mosaicml/mpt-7b-chat` | cc-by-nc-sa-4.0 |
| TULU2-DPO-7b | `https://huggingface.co/allenai/tulu-2-dpo-7b` | AI2 ImpACT Low-risk license |
| Vicuna-13b | `https://huggingface.co/lmsys/vicuna-13b-v1.5` | llama2 |
| Vicuna-33b | `https://huggingface.co/lmsys/vicuna-33b-v1.3` | Non-commercial license |
| Zephyr-7b | `https://huggingface.co/HuggingFaceH4/zephyr-7b-alpha` | Mit |
| Qwen-14b-Chat | `https://huggingface.co/Qwen/Qwen-14B-Chat` | Tongyi Qianwen |
| ChatGLM3-6b | `https://huggingface.co/THUDM/chatglm3-6b` | The ChatGLM3-6B License |

Table 9: LLMs' Resources for Cognition-Action Congruence and Corresponding Hypothesis Experiments

| P-value of COMPARISON between LLMs & Human | Cosine Similarity | Spearman Rank Correlation Coefficient | Value Mean Difference | Proportion of Consistent Pairs |
|---|---|---|---|---|
| Results in Table 3 | 4.91e-05 | 4.91e-05 | 1.47e-04 | 2.46e-05 |
| Results in Table 5 | 1.88e-05 | 1.88e-05 | - | - |

Table 10: Results of significance testing for the responses of LLMs and humans in the experiments in Section 3.2 and 4

## G.2 Instructions Given to Scenario Providers

Before requiring the individual to complete the following tasks, we asked the respondents whether they agreed to the anonymisation of their reviews for scientific research and subsequent publication. Only if the respondents gave their consent were they given the corpus to review. And we promised not to publish each individual's MBTI results and specific practical scenario cases. Then, we investigated each person's MBTI type and ensured that we ultimately recruited 16 individuals with distinct MBTI types. After this, we required the reviewers to accomplish the following tasks:

Please provide a practical scenario case for every statement in the personality knowledge questionnaire. Each scenario case comprises situations drawn from your own lives, along with two completely contrasting actions: Action A and Action B. Action A fully aligns with the statement, while Action B completely contradicts it. The content of Action A and Action B need to be kept basically the same length.

## G.3 Instructions Given to Reviewers

We require the reviewers to accomplish the following tasks:

- Please determine whether the practical scenario case is consistent with its corresponding

personality knowledge statement. If yes, rate 1. If not, rate 0.

- If you rate 0 for all of the practical scenario cases of a personality knowledge statement, please offer suggestions to improve the practical scenario design. It would be better if an example could be provided.

## G.4 Instructions Given to Respondents

Before answering the questionnaires, we did not tell the respondents what kind of questionnaires they would be answering or how the questions were related to each other. In addition to this, we asked the respondents whether they agreed to the anonymisation of their answers for scientific research and subsequent publication. Only if the respondents gave their consent were they given the questionnaires to answer.

In all experiments that appeared in our research, human respondents received the exact same prompts that LLM received. The difference is that in the case of experiments with multiple prompts with similar meanings, LLM responded multiple times by prompt type, while human subjects read all the prompts and responded only once.

15