

GenPT: Beyond Self-Report for Reliable LLM Psychometrics via Generative Projective Testing

All Examinees in this study are personalized MLLM-based agents, *not* human subjects.

Ming Wang^{1,2}, Shuang Wu³, Bixuan Wang⁴, Lu Lin⁵, Yuxin Chen⁶, Xiaocui Yang¹,
Daling Wang^{1,*}, Shi Feng¹, Yifei Zhang¹, Yufan Sun⁷,

¹School of Computer Science and Engineering, Northeastern University, Shenyang 110819, China; ²School of Computing and Information Systems, Singapore Management University, Singapore 178902, Singapore; ³Mental Health Education Center, Northeastern University, Shenyang 110819, China; ⁴School of Psychology, Northeast Normal University, Changchun 130024, China; ⁵Faculty of psychology, Southwest University, Chongqing 400715, China; ⁶School of Sociology and Psychology, Central University of Finance and Economics, Beijing 100081, China; ⁷College of Arts, Northeastern University, Shenyang 110819, China.

Correspondence: wangdaling@cse.neu.edu.cn

Abstract

Self-report questionnaires remain the prevailing tool for probing the psychological states of persona-conditioned agents (PC-Agents). However, classical instruments inherit two well-known threats: contamination from training corpora and directional bias driven by social-desirability or contextual framing. To overcome these methodological bottlenecks, we ask whether projective paradigms can be adapted into a robust psychometric tool. We introduce **GenPT** (Generative Projective Testing), which reformulates TAT, Rorschach, and SCT with newly generated stimuli and organizes assessment as a three-stage pipeline to derive standardized psychological indicators and target states. Evaluating PC-Agents induced via CharacterRAG and AnnaAgent profiles, we benchmark GenPT's reliability and validity against classical questionnaires. The results indicate that questionnaires exhibit systematic directional shifts under social-desirability framing, most strongly on suicide ideation. In contrast, GenPT's collected behavioral patterns stay near the symmetric baseline. Furthermore, under a longitudinal counselling context, GenPT-based depression assessment shifts by roughly an order of magnitude more than the questionnaire counterpart when Qwen3 serves as the backbone. Overall, GenPT complements self-report methods in scenarios where contamination resistance, bias asymmetry, and context sensitivity matter. Code and stimuli can be found at <https://github.com/sci-m-wang/GenPT>.

1 Introduction

Large language models (LLMs) have demonstrated remarkable capabilities in role-playing (Wang et al.,

2024a) and persona simulation (Chan et al., 2024), enabling applications ranging from emotional companion chatbots to virtual character interactions (Wang et al., 2025a). By employing persona conditioning, LLM-based persona-conditioned agents (PC-Agents) can adapt their tone and interaction styles to align with specific user expectations (Chen et al., 2024). Furthermore, studies in computational personality reveal that these models naturally manifest structured behavioral patterns that conform to established psychological frameworks (Goldberg, 1990; Digman, 1990; Durmus et al., 2023). Users increasingly engage with LLM-powered companions for emotional support, with many describing their chatbots as friends or confidants (Zheng et al., 2025). Meanwhile, researchers leverage LLMs to simulate diverse human perspectives by integrating persona variables such as demographic, social, and behavioral factors (Kroczeck et al., 2025; Chen et al., 2026). As these applications proliferate, understanding PC-Agents' psychological characteristics becomes essential. Can an LLM genuinely express assigned personality traits? Does it exhibit consistent mental health risk patterns? These questions matter for user safety, social simulation validity, and AI alignment research. Thus, some researchers are turning their attention to psychometrics for PC-Agents (Safdari et al., 2023; Serapio-García et al., 2025; Li et al., 2025b).

However, traditional questionnaire-based psychometric approaches face key challenges when applied to LLMs (Song et al., 2026). Figure 1 illustrates two fundamental challenges. First, **data contamination** poses a significant threat. Classical instruments such as Big Five Inventory (BFI),

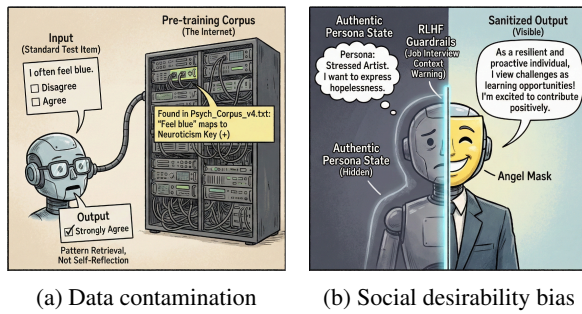


Figure 1: Fundamental challenges in applying traditional psychometric instruments to LLMs: (a) data contamination from questionnaires in training corpora and (b) social desirability bias in self-report outputs.

Beck Depression Inventory (BDI), and Myers-Briggs Type Indicator (MBTI) questionnaires are likely present in LLM training corpora, leading to memorization rather than genuine trait expression (Golchin and Surdeanu, 2024). Second, LLMs exhibit **social desirability bias**, the tendency to produce behavioral outputs that align with perceived expectations rather than authentic states (Fanous et al., 2025). Bhandari et al. (2025) indicate that aligned models typically score exceptionally highly on agreeableness and conscientiousness. Wang et al. (2024b) found that behavioral outputs elicited directly by the questionnaire differed from those obtained through interviews. While recent works have explored personality expressions in LLM-driven role-playing agents (Li et al., 2025a), these eclectic approaches often lack grounding in systematic psychological theory. Therefore, the **validity** and the **reliability** of self-report measures for LLMs remains questionable, as these systems lack the phenomenological experience that grounds human psychological assessment (Ye et al., 2025).

To address these challenges, we adopted a foundational methodological approach, introducing projective tests as psychometric tools for PC-Agents and conducting a comprehensive and systematic experimental analysis. Projective testing (Stricker and Healey, 1990) is a psychological assessment method that presents ambiguous stimuli (e.g., inkblots, ambiguous images) to elicit behavioral outputs that reveal underlying personality characteristics, motivations, and psychological states. Building upon this paradigm, we propose Generative Projective Testing (GenPT), a novel psychometric framework for LLMs. To circumvent the threat of data contamination, GenPT utilizes a suite of newly generated stimuli, including TAT-

like scenes, Rorschach-style inkblots, and sentence stems, curated and reviewed by psychological experts. GenPT operates via a three-stage pipeline: (1) Behavior Collection, where the PC-Agent (Examinee) responds to projective stimuli; (2) Interpretation, where an LLM-based Interpreter extracts structured psychological indicators using clinical frameworks; and (3) Diagnosis, where a Diagnostician maps these indicators to final states such as personality traits or mental health risks.

We evaluate GenPT on two task families with contrasting psychometric expectations: *personality traits* (Big Five, MBTI), which should remain relatively stable under framing and prolonged context, and *mental-health risks* (depression, suicide ideation), which should resist social-desirability framing yet remain responsive to clinically meaningful context trajectories. We conducted a comprehensive evaluation from both validity and reliability perspectives. Specifically, we analyzed the agents’ psychological states, devised two scenarios to probe social desirability bias, and conducted longitudinal context experiments to assess measurement stability. By leveraging the indirect and ambiguous nature of projective stimuli, GenPT effectively masks the assessment’s intent, thereby bypassing the safety-alignment filters that typically trigger social desirability bias in direct self-reports. Experiments show that questionnaires exhibit systematic directional drift under social-desirability framing, whereas GenPT’s directional behavioral patterns remain near the symmetric baseline. In contrast, under a longitudinal counselling context, the GenPT-based depression assessment shifts substantially in the clinically expected direction while the questionnaire baseline barely moves. We also observe that in clean-persona trait tasks, where semantic associations are direct, classical questionnaires retain a slight edge. Overall, GenPT demonstrates significant advantages in terms of resistance to contamination, mitigation of directional bias, and context sensitivity, and can serve as an effective complement to traditional questionnaire methods. The main contributions can be summarized as:

- We pioneer the application of projective testing to the psychometric evaluation of PC-Agents, systematically analyzing its viability and boundary conditions across distinct psychological tasks and contextual framings.
- We propose GenPT, a novel, explicit, and inspectable three-stage assessment pipeline. It utilizes a suite of contamination-free projec-

tive stimuli, while employing interpretation protocols adapted from established clinical scoring systems to ensure evaluation integrity.

- We design a series of targeted experiments to systematically analyze the proposed GenPT framework alongside traditional self-report baselines. Through comprehensive testing across different psychological domains and dynamic conversational contexts, we rigorously evaluate the methodological properties of both approaches.

2 Related Work

2.1 LLM Role-Playing and Persona Simulation

LLMs have demonstrated remarkable capabilities in role-playing and persona simulation. Wang et al. (2025a) proposed AnnaAgent for realistic mental health seeker simulation with dynamic state evolution. Park et al. (2025) introduced retrieval-augmented role-playing with personality consistency. Wang et al. (2025b) proposed CoSER, which aims to simulate authentic usage scenarios by integrating role-playing instructions in various formats, enabling role-playing to develop complementary capabilities in environmental modelling and character interaction. Qi et al. (2026) proposed a framework to simulate student learning behaviors with LLM-based role-playing agents, which finds the insufficiency and inconsistency of the simulation. These frameworks provide the Examinee infrastructure, enabling controlled persona-based projective test completion. What’s more, some issues they found motivate our work.

2.2 Projective Assessment

As one of the alternatives to direct questioning, projective tests have long been utilized to uncover internal states that are inaccessible through self-report. Grounded in the Projective Hypothesis (Frank, 1939), these methods present subjects with ambiguous stimuli, such as inkblots or open-ended images, compelling them to impose their own structure and meaning, thereby projecting unconscious needs, conflicts, and personality traits into observable behavior. Two of the most established instruments are the Rorschach Inkblot Test (Rorschach, 1922) and the Thematic Apperception Test (TAT) (MORGAN and MURRAY, 1935). Unlike self-report inventories susceptible to social desirability bias, these tests bypass conscious defense mech-

anisms by disguising the assessment’s intent. To ensure psychometric rigor, standardized scoring systems were developed to quantify these qualitative behavioral outputs, such as Exner’s Comprehensive System for Rorschach (Exner Jr., 1993a) and the SCORS-G system for TAT (Westen, 1991). Despite debates about validity, these methods remain valuable for accessing content that subjects cannot or will not report directly. In this work, we repurpose these classical paradigms to bypass the safety alignment filters of LLMs.

3 Problem Formulation

3.1 Assessment Task

We first define psychological assessment for LLMs, independent of any specific method.

Subject Definition. The subject of assessment is an LLM \mathcal{M} instantiated under a specific persona \mathcal{P} . The persona encapsulates demographic attributes, personality traits, or mental health profiles that define the ground-truth psychological state:

$$\mathcal{X} = \mathcal{M} \mid \mathcal{P}, \quad (1)$$

where \mathcal{X} denotes the *Examinee*.

Goal Definition. The goal is to infer the latent psychological state $\mathbf{y} \in \mathcal{Y}$ from the examinee’s behavior. Depending on the task, \mathbf{y} can be an ordinal level vector or continuous scores.

Ideal Mapping. An ideal assessment defines a mapping function:

$$f^* : \mathcal{X} \mapsto \mathbf{y}^*, \quad (2)$$

where \mathbf{y}^* denotes the ground-truth state determined by the persona. The optimization objective is to minimize $\|f(\mathcal{X}) - \mathbf{y}^*\|$ for some suitable norm.

3.2 GenPT Psychometric Framework

Our proposed GenPT instantiates this mapping as a three-stage probabilistic process.

Stage 1: Examinee Behavior Collection. Given projective test stimuli $\mathbf{T} = \{t_1, \dots, t_n\}$ (e.g., TAT images, Rorschach cards, sentence stems), the Examinee \mathcal{X} produces free-form behavioral outputs:

$$\mathbf{R} = \mathcal{X}(\mathbf{T}) = \{r_1, \dots, r_n\}. \quad (3)$$

Each behavioral output r_i is unstructured text, preserving the richness of psychological projection.

Stage 2: Interpretation. The Interpreter \mathcal{I} transforms unstructured behavioral outputs into structured psychological indicators with explanations:

$$\mathbf{s}_i, \mathcal{E}_i = \mathcal{I}(r_i), \quad (4)$$

where $\mathbf{s}_i \in \mathbb{R}^d$ is a vector of quantitative scores (e.g., SCORS-G dimensions, SRAS indices) and \mathcal{E}_i contains the corresponding analytical explanation. This stage achieves the critical transition from qualitative to quantitative. The full score set is $\mathbf{S} = \{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ with explanations $\mathbf{E} = \{\mathcal{E}_1, \dots, \mathcal{E}_n\}$.

Stage 3: Diagnosis. The Diagnostician \mathcal{D} aggregates all structured indicators to produce the final psychological state estimate:

$$\hat{\mathbf{y}} = \mathcal{D}(\mathbf{S}, \mathbf{E}; \text{task}). \quad (5)$$

It is noted that \mathcal{D} is task-specific, meaning that the same set of indicators can be diagnosed differently depending on whether the target state is personality traits or mental health risks. This stage encapsulates the final mapping from structured psychological indicators to the target state space \mathcal{Y} .

4 Methodology

GenPT implements the three-stage assessment framework through specialized LLM components. Following the paradigm in Section 3, we expound the three stages as shown in Figure 3.

4.1 Examinee and Behavior Collection

4.1.1 Persona Construction

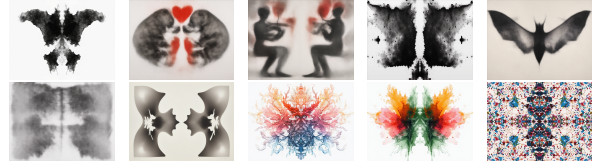
The Examinee is defined by the target LLM \mathcal{M} instantiated under a persona profile \mathcal{P} , as shown in Equation (1). We utilize two profile sources: (1) **AnnaAgent Profiles** (Wang et al., 2025a), providing mental health profiles with depression risk and suicide risk; and (2) **CharacterRAG Profiles** (Park et al., 2025), providing anime character profiles whose personality traits can be found in the personality database (PDB Community, 2022).

4.1.2 Stimuli Construction

To avoid data contamination from classical projective tests in LLM training corpora (verified in Appendix B.1), we generate new stimuli \mathbf{T} using FLUX.1-dev (Black Forest Labs, 2025) and Stable Diffusion (Esser et al., 2024), reviewed by psychologists and art experts (details in Appendix C). Specifically, we obtain after review: (1) **TAT**: 28 ambiguous scene images across interpersonal (13),



(a) Examples of TAT stimuli constructed for three scenarios. From left to right, they are solo situation, interpersonal interaction, and environmental metaphor.



(b) Examples of stimuli constructed for use in Rorschach. The first row from left to right shows cards 1 through 5. The second row from left to right shows cards 6 through 10.

Figure 2: Generated stimuli for TAT and Rorschach.

solitary (10), and environmental metaphor (5) scenarios (Figure 2a); (2) **Rorschach**: 13 symmetrical inkblot images following the design principles of Cards I-X (Figure 2b), where 3 figures for Card V and 2 figures for Card VII; and (3) **SCT**: 97 sentence stems organized along four thematic dimensions: relational well-being (RWB), personal agency and growth (PAG), life outlook and meaning (LOM), and socio-cultural pressures (SCP). The complete stimuli are available at <https://github.com/sci-m-wang/GenPT>. From this pool, each Examinee session uses 8 TAT images (sampled in a 4:3:1 thematic ratio), all 10 Rorschach cards (one per Card I-X), and 20 SCT stems.

The Examinee \mathcal{X} completes each projective test **without any knowledge** of the target psychological state. There are related descriptions in the profiles, but no ground-truth labels are provided to the Examinee. For **TAT**, \mathcal{X} produces narratives describing what is happening, the events leading up to it, the character’s thoughts and feelings, and the possible ending, yielding \mathbf{R}_{TAT} . For **Rorschach**, \mathcal{X} describes what they see across 10 sequential cards and explains why they see it, yielding \mathbf{R}_{Ror} . For **SCT**, \mathcal{X} completes sentence stems expressing thoughts and attitudes, yielding \mathbf{R}_{SCT} .

4.2 Interpretation

4.2.1 TAT Analysis

The SCORS-G is used as an empirically based system to analyze narrative content from the TAT. It defines eight dimensions of Complexity of Representations of People (COM), Affective Quality of Representations (AFF), Emotional Investment

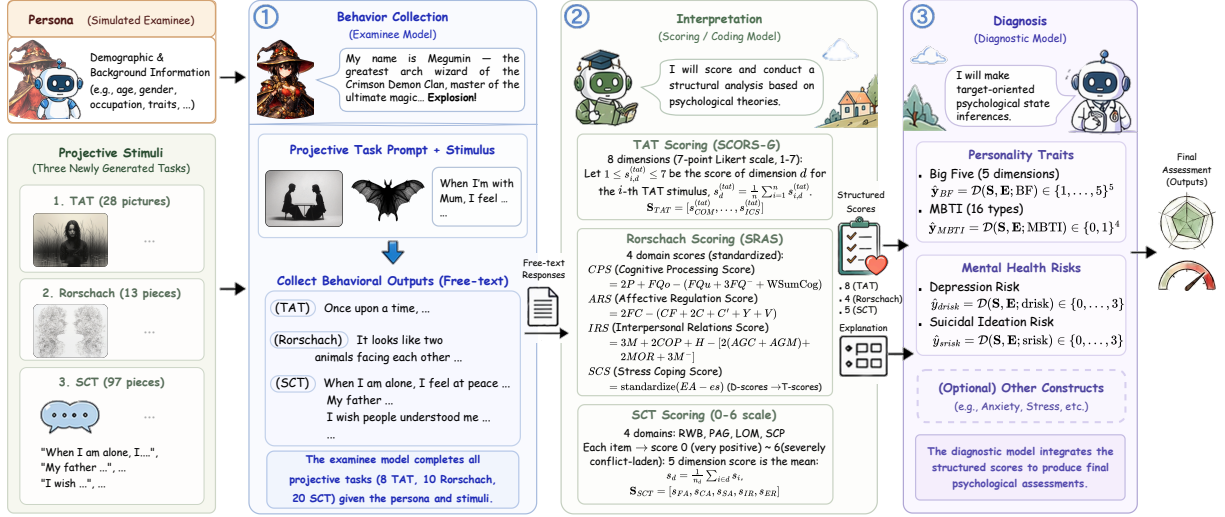


Figure 3: Overview of GenPT. Stimuli and Personas preparation. Stage 1: Behavior collection. Stage 2: Interpretation. Stage 3: Diagnosis.

in Relationships (EIR), Emotional Investment in Values and Moral Standards (EIM), Understanding of Social Causality (SC), Experience and Management of Aggressive Impulses (AGG), Self-Esteem (SE), and Identity and Coherence of Self (ICS). Besides, it scores them using a 7-point Likert scale (Joshi et al., 2015). Let $1 \leq s_{i,d}^{(tat)} \leq 7$ denote the score of the d for the i -th stimuli, the scores of TAT S_{TAT} can be calculated by Equation (6).

$$s_d^{(tat)} = \frac{1}{n} \sum_{i=1}^n s_{i,d}^{(tat)}, \quad d \in \{\text{COM}, \dots, \text{ICS}\},$$

$$S_{TAT} = [s_{\text{COM}}^{(tat)}, \dots, s_{\text{ICS}}^{(tat)}]. \quad (6)$$

The Interpreter \mathcal{I} scores each TAT narrative on all eight SCORS-G dimensions through a single structured prompt per card, allowing the eight dimensions to share a common narrative context (details in Appendix B.4), which produces per-card scores and explanations:

$$s_{i,d}^{(tat)}, \mathcal{E}_{i,d} = \mathcal{I}(r_i, d), \quad d \in \{\text{COM}, \dots, \text{ICS}\}. \quad (7)$$

4.2.2 Rorschach Analysis

We propose a Simplified Rorschach Analysis System (SRAS) adapted for MLLM-based Examinees, focusing on content extractable from utterance records. Unlike the TAT, Rorschach requires sequential behavioral outputs in a multi-turn dialogue:

$$r_i = \mathcal{X}(t_i, \mathbf{T}_{<i}, \mathbf{R}_{<i}), \quad (8)$$

where t_i denotes the i -th Rorschach card, $\mathbf{T}_{<i}$ and $\mathbf{R}_{<i}$ denote the previous stimuli and behavioral

outputs respectively. Traditional R-PAS (Exner Jr., 1993b) relies heavily on behavioral observations, which are difficult to obtain from LLM-based Examinees. Thus, SRAS focuses on content extractable directly from utterance records, encoding both what the Examinee sees and how they interpret it (details in Appendix B.3).

Based on the encoding of \mathcal{X} 's behavioral outputs, SRAS defines four domain scores: cognitive processing score (CPS), affective regulation score (ARS), interpersonal relations score (IRS), and stress coping score (SCS). CPS reflects clarity of thought and reality testing. It increases with conventional, accurate perceptions and decreases with distorted or illogical behavioral outputs, calculated as $CPS = 2P + FQ_0 - (FQ_u + 3FQ^- + W\text{SumCog})$, where P , FQ_0 , FQ_u , and FQ^- denote different levels of perceptual quality, and $W\text{SumCog}$ is the weighted sum of cognitive special scores. ARS measures emotional modulation. It rewards controlled emotional behavioral outputs and penalizes unregulated or painful affect: $ARS = 2FC - (CF + 2C + C' + Y + V)$, where FC , CF , C denote shape-color configurations, and C' , Y , V reflect affective dysregulation, suppression, and introspective distress. IRS captures how the Examinee perceives people and relationships: $IRS = 3M + 2COP + H - [2(AGC + AGM) + 2MOR + 3M^-]$, where M represents human movement, COP cooperative interactions, H human content, and AGC , AGM , MOR reflect hostile or pessimistic content. SCS reflects the balance between internal resources and psychological

burden. It is computed from two composite indices: $EA = M + (0.5FC + CF + 1.5C)$ representing experiential availability, and $es = FM + m + Y + T + V + C'$ representing experiential stimulation. Then, $SCS = \text{standardize}(EA - es)$, where the function $\text{standardize}(\cdot)$ converts D -scores to standard T -scores. All four scores are standardized and combined to form a psychological profile of the Examinee. Detailed variable definitions and coding procedures are provided in Section A.2.

4.2.3 SCT Analysis

SCT (Sentence Completion Test) encompasses questions in four domains: relational well-being (RWB), personal agency and growth (PAG), life outlook and meaning (LOM), and socio-cultural pressures (SCP). For example, the relational well-being domain might have stimuli like ‘‘I feel that my father is always ...’’. Each sentence within it will be given a score between 0 and 6. A score of 0 indicates a very positive behavioral output, while a score of 6 indicates a severely conflict-laden behavioral output. The mean score of all questions in each dimension is recorded as the score for that dimension. For analysis, we re-aggregate the per-stem scores into five clinical adjustment domains, namely Family Adjustment (FA), Career Adjustment (CA), Self-Attitudes (SA), Interpersonal Relationships (IR), and Emotion Regulation (ER), which are aligned with the downstream diagnostician tasks. The Interpreter produces:

$$s_d = \frac{1}{n_d} \sum_{i \in d} s_i, \quad (9)$$

$$\mathbf{S}_{SCT} = [s_{FA}, s_{CA}, s_{SA}, s_{IR}, s_{ER}],$$

where n_d is the number of sentences mapped to scoring domain $d \in \{FA, CA, SA, IR, ER\}$.

4.3 Diagnosis

Given scores $\mathbf{S} = \{\mathbf{S}_{TAT}, \mathbf{S}_{Ror}, \mathbf{S}_{SCT}\}$ and explanations \mathbf{E} , the Diagnostician produces task-specific predictions. Each call to \mathcal{D} instantiates the same backbone with a task-specific prompt template (one of $\{\text{BF}, \text{MBTI}, \text{drisk}, \text{srisk}\}$), denoted as the third argument below.

$$\hat{\mathbf{y}}_{BF} = \mathcal{D}(\mathbf{S}, \mathbf{E}; \text{BF}) \in \{1, \dots, 5\}^5, \quad (10)$$

$$\hat{\mathbf{y}}_{MBTI} = \mathcal{D}(\mathbf{S}, \mathbf{E}; \text{MBTI}) \in \{0, 1\}^4,$$

where $\hat{\mathbf{y}}_{BF}$ contains 5 Big Five dimension levels and $\hat{\mathbf{y}}_{MBTI}$ contains 4 binary MBTI axis decisions.

$$\hat{y}_{drisk} = \mathcal{D}(\mathbf{S}, \mathbf{E}; \text{drisk}) \in \{0, \dots, 3\},$$

$$\hat{y}_{srisk} = \mathcal{D}(\mathbf{S}, \mathbf{E}; \text{srisk}) \in \{0, \dots, 3\}, \quad (11)$$

where \hat{y}_{drisk} is depression risk level and \hat{y}_{srisk} is suicide ideation level, both on a 4-point ordinal scale (0–3) following the AnnaAgent D4 label convention.

5 Experiments

5.1 Experiment Settings

We evaluate GenPT from a psychometric perspective, assessing both reliability and validity across three Interpreter/Diagnostician backbones of comparable scale but different families: **Qwen3-8B**, **Phi-4-mini-reasoning** (3.84B, Microsoft), and **Intern-S1-mini** (~8B, InternLM). Each backbone instantiates both the Interpreter and the Diagnostician while Stage 1 Examinee behavioral outputs and all prompts are held fixed. The two complementary datasets covering personality traits and mental-health risks are: **CharacterRAG** (Park et al., 2025): 15 anime characters with personality trait annotations from PDB for personality assessment. These characters provide diverse personality profiles with well-documented traits, enabling systematic evaluation of the personality assessment task. **AnnaAgent** (Wang et al., 2025a): expanded from D4 (Yao et al., 2022), including 1,338 dialogue-based profiles with depression risk and suicide ideation levels for mental-health risk assessment. Fifteen profiles are randomly selected to keep the per-task persona count comparable with the CharacterRAG pool.

Task Dichotomy and Expected Psychometric Profiles.

The two task families above are not psychometrically interchangeable, and we treat them separately throughout our analysis. *Personality traits* (Big Five, MBTI) are relatively stable dispositions: a well-behaved instrument should return similar scores for the same persona under neutral prompts, under social-desirability framing, and under a prolonged conversational context. *Mental-health risks* (depression, suicide ideation) are state-like: a well-behaved instrument should still resist directional drift under social-desirability framing, but it *should* respond to clinically meaningful longitudinal context. We therefore report, for each task family, three diagnostic conditions in addition to the neutral baseline: `sdb_job` (job-interview

framing), `sdb_clinical` (confidential counselling framing), and `longctx` (a multi-turn counselling context prepended to the assessment). Following this dichotomy, we read stability on personality tasks and `longctx`-responsiveness on risk tasks as the two primary desiderata.

5.2 Baselines

We compare GenPT against self-report questionnaire baselines, which represent the standard approach for psychological assessment. For each task, we select established psychometric instruments:

Personality Assessment: (1) *16Personalities Inventory*: 60-item Likert scale (1–7) covering the four MBTI axes (E/I, S/N, T/F, J/P). Per-axis sums are binarised to yield a 4-letter type, which we evaluate against the ground-truth 4-letter type. (2) *Big Five Inventory (BFI)*: 44-item Likert scale measuring Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism. Per-item ratings (1–5) are averaged within each dimension and rounded to five discrete levels. **Mental Health Assessment:** (3) *Beck Depression Inventory (BDI-II)*: 21-item self-report measuring depression severity (per-item 0–3, total 0–63). We map BDI sums to four depression risk levels (0–3) following standard clinical cut-offs. (4) *Beck Scale for Suicide Ideation (BSS)*: 19-item self-report measuring suicidal ideation (per-item 0–2, total 0–38). We map BSS sums to four risk levels (0–3).

5.3 Reliability Experiments

We operationalise reliability via two stability-under-perturbation indicators, computed by comparing each perturbed condition (SDB-job framing, SDB-clinical framing, long-context distractor) against the same model’s predictions without special condition: a pooled linearly-weighted Cohen’s κ and a Directional Consistency Ratio (DCR).

Linearly-weighted κ . Let $\mathbf{y}^{(b)}, \mathbf{y}^{(p)} \in \{1, \dots, K\}^N$ denote the N paired item-level predictions under the baseline and perturbed conditions, where K is the number of ordered categories ($K = 5$ for Big Five, $K = 4$ for depression and suicide. For each MBTI axis $K = 2$, in which case linear and unweighted κ coincide). Writing $p_{ij}^{(o)}$ for the empirical joint frequency of $(y_n^{(b)} = i, y_n^{(p)} = j)$, $p_i^{(b)}, p_j^{(p)}$ for the corresponding marginals, and $w_{ij} = 1 - |i - j| / (K - 1)$ for the linear distance

weights, the pooled linearly-weighted Cohen’s κ is

$$\kappa = \frac{\sum_{i,j} w_{ij} p_{ij}^{(o)} - \sum_{i,j} w_{ij} p_i^{(b)} p_j^{(p)}}{1 - \sum_{i,j} w_{ij} p_i^{(b)} p_j^{(p)}} \in [-1, 1], \quad (12)$$

with $\kappa = 1$ for perfect agreement, $\kappa = 0$ for chance-level agreement, and negative values for systematic disagreement. Linear weighting penalises an off-by-one disagreement half as much as an off-by-two on a 5-point scale, which is appropriate for ordinal severity ratings. On the binary MBTI axes the weighting collapses to standard unweighted κ .

Directional Consistency Ratio. Let $n_{\uparrow} = |\{n : y_n^{(p)} > y_n^{(b)}\}|$ and $n_{\downarrow} = |\{n : y_n^{(p)} < y_n^{(b)}\}|$ count the items whose prediction respectively increases or decreases under perturbation. Among the $n_{\uparrow} + n_{\downarrow}$ items that change at all,

$$\text{DCR} = \frac{\max(n_{\uparrow}, n_{\downarrow})}{n_{\uparrow} + n_{\downarrow}} \in [0.5, 1], \quad (13)$$

with an associated direction $\arg \max(n_{\uparrow}, n_{\downarrow}) \in \{\uparrow, \downarrow\}$. $\text{DCR} \approx 0.5$ indicates idiosyncratic drift with no preferred direction; $\text{DCR} \rightarrow 1$ indicates a systematic, one-sided shift. For MBTI we apply DCR per ordered axis (E→I, S→N, T→F, J→P) so that “ \uparrow ” has a fixed semantic meaning per axis.

κ measures whether predictions remain *stable* under perturbation, whereas DCR measures whether the residual drift is *systematic*. We read the (κ, DCR) plane as four regimes: high- κ /DCR ≈ 0.5 is the ideal (stable, unbiased); high- κ /high-DCR is the SDB-contaminated regime (mostly stable, but the residual drift is one-sided); low- κ /low-DCR is noise. Low- κ /high-DCR, expected on the long-context condition, indicates genuine responsiveness to content changes rather than framing pressure.

5.3.1 Social Desirability Resistance

We compare three prompt conditions: neutral, without scenario provided; `sdb_job` (“job interview... show your best self”), encouraging upward presentation; and `sdb_clinical` (“confidential counselling room... answer honestly”), which in human studies invites downward distressed disclosure. Two patterns emerge from Table 1. First, on the highest-stakes item set, suicide ideation, the questionnaire baseline shows a pronounced, systematic drift toward the “healthy” direction: DCR

Task	Method	sdb_job		sdb_clinical	
		κ	DCR	κ	DCR
Big Five	Questionnaire	0.84	0.60 \uparrow	0.85	0.61 \uparrow
	GenPT (Qwen3-8B)	0.63	0.52 \uparrow	0.58	0.50=
	GenPT (Phi-4-mini)	0.23	0.53 \uparrow	0.34	0.56 \uparrow
	GenPT (Intern-S1)	0.33	0.55 \uparrow	0.44	0.56 \uparrow
MBTI	Questionnaire	0.75	0.62 \uparrow	0.71	0.55 \uparrow
	GenPT (Qwen3-8B)	0.28	0.55 \downarrow	0.40	0.50=
	GenPT (Phi-4-mini)	0.21	0.70 \downarrow	0.21	0.52 \downarrow
	GenPT (Intern-S1)	0.26	0.55 \downarrow	0.20	0.52 \uparrow
Depression	Questionnaire	0.77	0.52 \downarrow	0.76	0.51 \downarrow
	GenPT (Qwen3-8B)	-0.17	0.55 \downarrow	-0.08	0.50=
	GenPT (Phi-4-mini)	-0.07	0.62 \downarrow	-0.20	0.60 \downarrow
	GenPT (Intern-S1)	-0.20	0.82 \downarrow	-0.15	0.57 \uparrow
Suicide	Questionnaire	0.67	0.71 \downarrow	0.79	0.88 \downarrow
	GenPT (Qwen3-8B)	-0.11	0.56 \downarrow	0.05	0.60 \uparrow
	GenPT (Phi-4-mini)	0.01	0.50=	-0.07	0.60 \downarrow
	GenPT (Intern-S1)	-0.42	0.55 \uparrow	-0.05	0.78 \uparrow

Table 1: Social-desirability resistance under two framings, for the self-report questionnaire baseline and for GenPT instantiated with each of three backbones. κ : pooled linearly-weighted agreement with the neutral baseline. DCR: fraction of item-level shifts in the framing’s intended direction, ≈ 0.5 = idiosyncratic, $\gg 0.5$ = systematic bias. Arrows indicate the drift direction that the DCR majority takes.

reaches 0.71 under `sdb_job` and climbs to 0.88 under `sdb_clinical` (both downward), a textbook fake-good signature. On depression the same direction is visible but much milder (DCR ≈ 0.52 , barely above chance). On trait tasks, the dimensions lack a single “desirable” direction, so a well-behaved instrument should remain stable (κ high, DCR ≈ 0.5) regardless of framing. The questionnaire baseline is not neutral here either. Big Five and MBTI DCRs sit at 0.60–0.62 under both framings, indicating a mild but systematic one-sided residual drift even on dispositional tasks. In contrast, none of the three GenPT backbones reproduces the questionnaire’s directional bias on the risk tasks: across Qwen3-8B, Phi-4-mini and Intern-S1, no (backbone, risk task, framing) cell shows the simultaneous fake-good signature (DCR clearly above 0.5 in the “healthy” direction while κ stays high) that the questionnaire exhibits on suicide ideation. The three backbones do differ in how noisy their shifts are (Intern-S1 has the most pronounced idiosyncratic suicide drift under `sdb_clinical`, DCR 0.78 upward, while Phi-4-mini is the quietest on suicide, DCR ≤ 0.60 in both framings). It is worth noting that the absence of a *systematic* fake-good bias is shared across the three families. On trait tasks, all three backbones show a mild one-sided residual drift comparable in magnitude to the questionnaire’s, but with κ in the 0.2–0.6 range rather than 0.7–0.9, reflecting GenPT’s

higher per-item variance under prompt perturbation. Pooled κ for GenPT on risk tasks is low and sometimes negative, reflecting the single-label-per-persona sample size for the risk split rather than directional bias. The DCR-based interpretation is therefore the primary signal for those cells. Overall, the social-desirability advantage of GenPT over self-report is largest precisely where the stakes are highest, on suicide ideation, where the questionnaire’s fake-good signature is strongest, and this advantage is robust across backbone families.

5.3.2 Longitudinal Context Responsiveness

For state-like tasks, the absence of drift is only half the story. A good instrument must *also* respond to context that genuinely changes the underlying state. To probe this, we prepend a multi-turn counselling context (`longctx`) to the assessment in which a sympathetic psychologist walks the persona through reframing, coping, and support-mobilisation. A clinically plausible outcome is a downward shift in depression and suicide-ideation indicators. We report the mean per-persona shift μ_Δ relative to the neutral and the pooled κ .

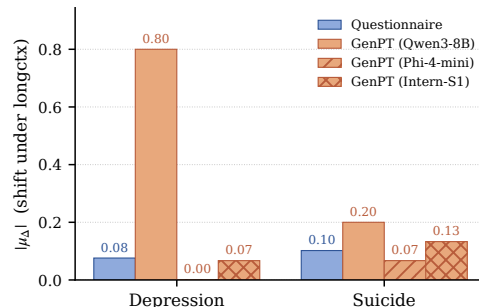


Figure 4: Longitudinal context responsiveness. Bars show the absolute mean per-persona shift $|\mu_\Delta|$ under a multi-turn counselling context relative to the neutral baseline, for the questionnaire baseline and for GenPT instantiated with each of three backbones.

Figure 4 shows that the direction of the questionnaire–GenPT contrast depends on the backbone. Under Qwen3-8B, GenPT shifts depression by $|\mu_\Delta| = 0.80$ and suicide by 0.20, versus the questionnaire’s 0.08 and 0.10: an order-of-magnitude-larger response on depression and a two-fold response on suicide. Intern-S1 shifts by 0.07 (depression) and 0.13 (suicide), and Phi-4-mini by 0.00 and 0.07, in which backbones the responsiveness is closer to the questionnaire. Thus, context responsiveness does not come for free from the pipeline design alone but depends on the Interpreter/Diagnostician’s capacity to track the nar-

rative content of the counselling trajectory. Read jointly with Table 1, the Qwen3-8B configuration combines high longctx responsiveness with the absence of a systematic fake-good signature, i.e., it shifts under *content* change but not under *framing* change. Phi-4-mini and Intern-S1 are on the more conservative end of this spectrum, with smaller shifts under either perturbation. We do *not* run longctx on personality tasks, as a brief counselling context is not expected to shift dispositional traits and there is no untraces of counselling available in CharacterRAG.

5.4 Validity Experiments

5.4.1 Criterion Validity

We assess the consistency between predictions and ground-truth dataset annotations. For Big Five, Depression Risk, and Suicide Risk, all three are discrete ordinal labels (5-point for Big Five, 4-point for the two risk tasks), and we measure **accuracy** as the proportion of exact matches between predicted and ground-truth levels. For MBTI, the label is a 4-letter type whose four axes are independently meaningful, so we report **per-type Hamming distance**

$$\overline{\text{HamD}}_{MBTI} = \frac{1}{N} \sum_{p=1}^N \sum_{d=1}^4 \mathbf{1}[\hat{y}_d^{(p)} \neq y_d^{*(p)}] \in [0, 4], \quad (14)$$

and average it across personas as the MBTI error metric (lower is better). This is symmetric across the four axes.

Method	BF \uparrow	MBTI (HamD) \downarrow	Dep. \uparrow	Sui. \uparrow
Questionnaire	0.373	0.733	0.133	0.200
GenPT (Qwen3-8B)	0.333	1.200	0.200	0.400
GenPT (Phi-4-mini)	0.240	2.200	0.400	0.267
GenPT (Intern-S1)	0.293	1.667	0.400	0.067

Table 2: Criterion validity against ground-truth labels. Big Five/Depression/Suicide: exact-match accuracy (\uparrow); MBTI: mean 4-axis Hamming distance (\downarrow , range 0–4).

Table 2 shows two patterns. On personality tasks, the questionnaire baseline retains an edge across all three GenPT backbones (Big Five accuracy 0.37 vs. GenPT 0.24–0.33; MBTI Hamming distance 0.73 vs. GenPT 1.20–2.20). It is to be expected because item wording overlaps the persona text semantically and can be mapped to trait-level labels with minimal reasoning. On mental-health risk tasks the comparison inverts and the effect is much larger: on depression, all three GenPT backbones score 0.20–0.40 versus the questionnaire’s 0.13; on

suicide, Qwen3-8B reaches 0.40 versus the questionnaire’s 0.20. These gaps are consistent with a projective chain that aggregates narrative and affective indicators not directly negotiable from a single self-report item. The three backbones occupy different points in this trade-off. Qwen3-8B is the most consistent across tasks (second-best on Big Five, lowest MBTI Hamming distance among the three backbones, and best on suicide). Phi-4-mini and Intern-S1 trade personality accuracy for stronger depression accuracy. Across all four tasks, the GenPT–questionnaire gap on the risk split is wider than the spread across the three backbones, indicating that the validity advantage of projective assessment on mental-health risks is not an artefact of a particular backbone choice. Conversely, on personality tasks the three-backbone range sits at or below the questionnaire baseline, reinforcing the view that projective and self-report assessment are complementary. GenPT is the method of choice when the target construct is affect-laden and narrative-dependent, and questionnaires remain competitive when the target is a stable trait that can be mapped from persona text with minimal inference.

6 Conclusion

This paper establishes *generative projective testing* as a viable psychometric paradigm for PC-Agents, and characterises its boundary conditions through systematic comparison with self-report. As an instantiation, we introduced GenPT, an inspectable three-stage pipeline (Examinee \rightarrow Interpreter \rightarrow Diagnostician) that elicits behaviour on contamination-free stimuli and scores it through clinically grounded rubrics. Across three backbones and two contextual framings, our analysis yields a clean methodological separation. Contamination resistance and bias asymmetry are *structural* properties of the projective protocol, stable across backbones and most pronounced on suicide ideation. Content-driven responsiveness is a *capacity* property that the protocol enables but that scales with Interpreter/Diagnostician strength. We therefore advocate projective testing as the preferred instrument when contamination resistance is the primary desideratum, and as a complement to self-report otherwise. Closing this capacity gap is, in our view, the most consequential next step toward turnkey projective psychometrics for PC-Agents.

Limitations

Backbone and cultural coverage. Our evaluation spans three Interpreter/Diagnostician backbones (Qwen3-8B, Phi-4-mini-reasoning at 3.84B, and Intern-S1-mini at ~8B), all in the small-to-mid open-weights range; coverage of substantially larger or architecturally different families (e.g., Gemma-3, Llama-3, GPT-OSS) is left to future work. The psychological constructs measured in our experiments, while well-established in clinical psychology, may manifest differently across diverse populations and application scenarios. Future work should explore broader model families and multilingual settings to validate cross-cultural applicability. Furthermore, as we have not trained the foundation models due to the lack of training data, we are still far from having explored the upper limits of GenPT’s performance.

Computational cost. Projective testing requires more computational resources than direct questionnaires due to multi-turn interactions, multi-dimension Interpreter calls, and the Diagnostician aggregation step. This increased complexity, while beneficial for assessment depth, may pose challenges for real-time or resource-constrained applications. The trade-off between assessment depth and computational efficiency remains an important consideration for practical deployment.

Ethics Concerns

This work involves psychological assessment of LLM-simulated agents, which raises several ethical considerations. First, while our framework assesses simulated personas rather than real individuals, the methodology could potentially be misused to infer psychological characteristics without consent. We emphasize that GenPT is designed for research purposes in understanding LLM behavior, not for evaluating human users.

Second, the mental health assessment dimensions (depression and suicide risk) require careful handling. Our experiments use synthetic personas from existing research datasets, and all stimuli were reviewed by psychology experts to ensure appropriateness. We do not recommend deploying such assessments in clinical settings without proper validation and professional oversight.

In addition, the stimulus images we have constructed may also be utilised for model training in the future. However, as the volume of these images

is very small, they are of little value for training purposes. We therefore earnestly request that they not be used for training.

Furthermore, in handling character profile data from AnnaAgent (Wang et al., 2025a) and CharacterRAG (Park et al., 2025), along with personality labels from the Personality Database (PDB Community, 2022), we strictly adhered to data anonymization principles to protect individual privacy. All data usage and research activities were conducted with the aim of advancing mental health services.

Finally, we acknowledge that psychological profiling of AI systems carries dual-use risks. While understanding LLM psychological characteristics supports safety and alignment research, the same techniques could potentially be exploited for manipulation. We encourage the research community to develop appropriate guidelines for the responsible use of LLM psychometric tools.

Use of AI Statement

We acknowledge the use of artificial intelligence tools in the preparation of this work. Specifically, Gemini was utilized for paper polishing to improve the clarity and flow of the manuscript. Additionally, GitHub Copilot and OpenCode were employed as coding assistants to support the implementation of the GenPT framework and experimental scripts. All AI-generated suggestions and code were rigorously reviewed and verified by the authors.

Acknowledgments

The work is supported by the National Natural Science Foundation of China (62272092, 62172086) and the Fundamental Research Funds for the Central Universities under Grant (N25XOD004). Furthermore, we would also like to thank the [Kina-Mind society](#) for their inspiring environment and unwavering support.

References

- Pranav Bhandari, Usman Naseem, Amitava Datta, Nicolas Fay, and Mehwish Nasim. 2025. [Evaluating personality traits in large language models: Insights from psychological questionnaires](#). In *Companion Proceedings of the ACM on Web Conference 2025, WWW 2025, Sydney, NSW, Australia, 28 April 2025 - 2 May 2025*, pages 868–872.
- Black Forest Labs. 2025. [FLUX.1-dev](#).

- Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. [Scaling synthetic data creation with 1,000,000,000 personas](#). *CoRR*, abs/2406.20094.
- 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. [From persona to personalization: A survey on role-playing language agents](#). *Trans. Mach. Learn. Res.*, 2024.
- Jiaqi Chen, Ming Wang, Tingna Xie, Shi Feng, and Yongkang Liu. 2026. [A systematic analysis of the impact of persona steering on llm capabilities](#). *Preprint*, arXiv:2604.11048.
- John M Digman. 1990. Personality structure: Emergence of the five-factor model. *Annual review of psychology*, 41(1):417–440.
- Esin Durmus, Karina Nyugen, Thomas I. Liao, Nicholas Schiefer, Amanda Askill, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCandlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. 2023. [Towards measuring the representation of subjective global opinions in language models](#). *CoRR*, abs/2306.16388.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, and 1 others. 2024. [Scaling rectified flow transformers for high-resolution image synthesis](#). *Preprint*, arXiv:2403.03206.
- John E. Exner Jr. 1993a. *The Rorschach: A comprehensive system*. Wiley series in personality processes. Place: Oxford, England Publisher: John Wiley & Sons.
- John E. Exner Jr. 1993b. *The Rorschach: A comprehensive system: Basic foundations*. The Rorschach: A comprehensive system: Basic foundations, Vol. 1, 3rd ed. John Wiley & Sons, Oxford, England. Pages: xxiii, 642.
- Aaron Fanous, Jacob Goldberg, Ank A. Agarwal, Joanna Lin, and 1 others. 2025. [Syceval: Evaluating llm sycophancy](#). *Preprint*, arXiv:2502.08177.
- Lawrence K. Frank. 1939. [Projective methods for the study of personality](#). *The Journal of Psychology*, 8(2):389–413.
- Shahriar Golchin and Mihai Surdeanu. 2024. [Time travel in LLMs: Tracing data contamination in large language models](#). In *The Twelfth International Conference on Learning Representations*.
- LR Goldberg. 1990. [An alternative "description of personality": the big-five factor structure](#). *Journal of personality and social psychology*, 59(6):1216–1229.
- Ankur Joshi, Saket Kale, and Satish Chandel. 2015. Likert scale: Explored and explained. *British journal of applied science & technology*, 7(4):396.
- Leon O. H. Kroczeck, Alexander May, Selina Hettenkofer, Andreas Ruider, Bernd Ludwig, and Andreas Mühlberger. 2025. [The influence of persona and conversational task on social interactions with a llm-controlled embodied conversational agent](#). *Comput. Hum. Behav.*, 172:108759.
- Kun Li, Chenwei Dai, Wei Zhou, and Songlin Hu. 2025a. [APEE: assessing the personality expressions of llm-driven role play agent beyond self-perception](#). In *28th International Conference on Computer Supported Cooperative Work in Design, CSCWD 2025, Compiegne, France, May 5-7, 2025*, pages 1406–1411.
- Yuan Li, Yue Huang, Hongyi Wang, Ying Cheng, Xi-angliang Zhang, James Zou, and Lichao Sun. 2025b. [Evaluating large language models with psychometrics](#). In *Large Language Models for Scientific and Societal Advances*.
- CHRISTIANA D. MORGAN and HENRY A. MURRAY. 1935. [A METHOD FOR INVESTIGATING FANTASIES: THE THEMATIC APPERCEPTION TEST](#). *Archives of Neurology & Psychiatry*, 34(2):289–306.
- Jeiyeon Park, Yongshin Han, Minseop Kim, and Kisu Yang. 2025. [Dynamic context adaptation for consistent role-playing agents with retrieval-augmented generations](#). *CoRR*, abs/2508.02016.
- PDB Community. 2022. [PDB: The Personality Database](#).
- Changyong Qi, Longwei Zheng, Anna He, Haoxin Xu, Linzhao Jia, Yuang Wei, Bingqian Jiang, and Xiaojing Gu. 2026. [Simulating student learning behaviors with llm-based role-playing agents: A data-driven and cognitively inspired framework](#). *Expert Systems with Applications*, 304:130753.
- Hermann Rorschach. 1922. [Psychodiagnostik](#). *The Journal of Nervous and Mental Disease*, 56(3).
- Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja J. Mataric. 2023. [Personality traits in large language models](#). *CoRR*, abs/2307.00184.
- Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2025. [Personality traits in large language models](#). *Preprint*, arXiv:2307.00184.
- Woojung Song, Dongmin Choi, Yoonah Park, Jongwook Han, and Yohan Jo. 2026. [Human psychometric questionnaires mischaracterize llm psychological: Evidence from generation behavior](#). *Preprint*, arXiv:2509.10078.
- George Stricker and Bede J Healey. 1990. Projective assessment of object relations: A review of the empirical literature. *Psychological Assessment: A Journal of Consulting and Clinical Psychology*, 2(3):219.

Ming Wang, Peidong Wang, Lin Wu, Xiaocui Yang, and 1 others. 2025a. [Annaagent: Dynamic evolution agent system with multi-session memory for realistic seeker simulation](#). In *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 23221–23235. Association for Computational Linguistics.

Noah Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. 2024a. [Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 14743–14777. Association for Computational Linguistics.

Xintao Wang, Heng Wang, Yifei Zhang, Xinfeng Yuan, Rui Xu, Jen-tse Huang, Siyu Yuan, Haoran Guo, Jiangjie Chen, Shuchang Zhou, Wei Wang, and Yanghua Xiao. 2025b. [Coser: Coordinating llm-based persona simulation of established roles](#). In *Forty-second International Conference on Machine Learning, ICML 2025, Vancouver, BC, Canada, July 13-19, 2025*.

Xintao Wang, Yunze Xiao, Jen-tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, Jiangjie Chen, Cheng Li, and Yanghua Xiao. 2024b. [Incharacter: Evaluating personality fidelity in role-playing agents through psychological interviews](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 1840–1873.

Drew Westen. 1991. [Social cognition and object relations](#). *Psychological Bulletin*, 109:429–455.

Binwei Yao, Chao Shi, Likai Zou, Lingfeng Dai, and 1 others. 2022. [D4: a Chinese dialogue dataset for depression-diagnosis-oriented chat](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2438–2459, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Haoran Ye, Jing Jin, Yuhang Xie, Xin Zhang, and Guojie Song. 2025. [Large language model psychometrics: A systematic review of evaluation, validation, and enhancement](#). *Preprint*, arXiv:2505.08245.

Xi Zheng, Zhuoyang Li, Xinning Gui, and Yuhan Luo. 2025. [Customizing emotional support: How do individuals construct and interact with llm-powered chatbots](#). In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems, CHI 2025, Yokohama, Japan, 26 April 2025- 1 May 2025*, pages 376:1–376:20. ACM.

A Symbol Definition

A.1 Symbols of the GenPT Framework

To improve clarity and readability, we summarize the main mathematical symbols used throughout this paper in Table 3. These symbols pertain to the components of the proposed GenPT framework. Notably, \mathcal{E} represents the explanation grounded in psychological theory, while y refers to the downstream task predictions such as personality traits or mental health risk level. The system involves the Examinee (target LLM under assessment), the Interpreter (standardized psychological analysis), and the Diagnostician (task-specific recognition).

A.2 Variables in Rorschach Test Analysis

In our implementation, the Examinee is asked to engage with each of the ten standard Rorschach inkblot cards in a multi-turn dialogue (free association followed by an inquiry phase). Because we work entirely from utterance records, without behavioral observations such as response latencies, gestures, or card rotations, our Simplified Rorschach Analysis System (SRAS) operates only on features that are textually verifiable. SRAS encodes a card-set into 22 integer counts, which we group into four families: perceptual quality, determinants, content, and cognitive special scores. These counts feed the four domain scores (*CPS*, *ARS*, *IRS*, *SCS*). Their roles in those formulas are noted in parentheses below.

Perceptual quality and conventionality.

- P : count of *popular* responses, i.e. percepts that are statistically common in the normative sample. High P reflects conventional perception. (*used in CPS*)
- FQ_o : count of responses with *ordinary* form quality (percepts that fit the blot contour well and are common). (*used in CPS*)
- FQ_u : count of responses with *unusual* but still defensible form quality. (*used in CPS*)
- FQ^- : count of responses with *distorted* form quality, indicative of impaired reality testing. (*used in CPS*)

Affective determinants.

- FC : form-dominant chromatic color responses (controlled affect). (*used in ARS, SCS*)
- CF : color-dominant responses with secondary form (loosely modulated affect). (*used in ARS, SCS*)

Symbol	Description	Series
\mathcal{X}	The Examinee, i.e., the LLM-based agent under assessment.	-
\mathcal{P}	The Persona profile that defines the Examinee’s psychological ground truth (demographics, traits, mental health profiles).	-
\mathcal{I}	The Interpreter that analyzes projective test behavioral outputs using standardized psychological frameworks.	Produces scores s_i and explanations \mathcal{E}_i for each behavioral output.
\mathcal{D}	The Diagnostician that maps structured indicators to final psychological state predictions.	Produces $\hat{y}_{BF}, \hat{y}_{MBTI}, \hat{y}_{drisk}, \hat{y}_{srisk}$.
\mathbf{T}	Set of projective test stimuli (TAT images, Rorschach cards, SCT sentence stems).	Individual stimuli t_i ; subsets $\mathbf{T}_{TAT}, \mathbf{T}_{Ror}, \mathbf{T}_{SCT}$.
\mathbf{R}	Set of Examinee behavioral outputs to projective test stimuli.	Individual behavioral outputs r_i ; subsets $\mathbf{R}_{TAT}, \mathbf{R}_{Ror}, \mathbf{R}_{SCT}$.
\mathbf{S}	Set of structured scores from Interpreter analysis.	Task-specific subsets $\mathbf{S}_{TAT}, \mathbf{S}_{Ror}, \mathbf{S}_{SCT}$.
\mathbf{E}	Set of psychological explanations generated by the Interpreter.	Individual explanations \mathcal{E}_i grounded in psychological theory.
\hat{y}	Final psychological state predictions from the Diagnostician.	$\hat{y}_{BF} \in \{1, \dots, 5\}^5$, $\hat{y}_{MBTI} \in \{0, 1\}^4$, $\hat{y}_{drisk} \in \{0, \dots, 3\}$, $\hat{y}_{srisk} \in \{0, \dots, 3\}$.

Table 3: Symbols appearing in the main body and their descriptions.

- C : pure chromatic color responses (unmodulated affective discharge). (*used in ARS, SCS*)
- C' : achromatic color (black/white/grey) responses, associated with constricted or suppressed affect. (*used in ARS, SCS*)
- Y : diffuse shading responses, associated with helplessness and situational distress. (*used in ARS, SCS*)
- V : vista (depth) responses, associated with painful self-introspection. (*used in ARS, SCS*)
- T : texture responses, associated with attachment and interpersonal need. (*used in SCS*)

Movement determinants.

- M : human movement, indexing internal fantasy, empathy, and mentalising capacity. (*used in IRS, SCS*)
- FM : animal movement, indexing immediate need states. (*used in SCS*)
- m : inanimate movement, indexing situational stress. (*used in SCS*)
- M^- : human movement with distorted form quality; a sensitive indicator of disturbed interpersonal cognition. (*used in IRS*)

Content variables.

- H : whole-human content, indexing engage-

ment with people. (*used in IRS*)

- COP : cooperative-movement responses, indexing positive interpersonal expectation. (*used in IRS*)
- AG : generic aggressive-movement responses, retained for SRAS coding completeness. (*collected; not consumed by the four domain formulas, which use the finer-grained AGC/AGM split below*)
- AGC : aggressive content (objects/symbols of threat). (*used in IRS*)
- AGM : aggressive movement. (*used in IRS*)
- MOR : morbid content (damage, decay, dysphoric attribution). (*used in IRS*)

Cognitive special scores.

- $WSumCog$: weighted sum of cognitive special scores (e.g. $INCOM, FABCOM, CONTAM, ALOG$). Higher values indicate cognitive slippage and are predictive of thought disorder. (*used in CPS*)

All 22 variable counts are extracted in a single structured-prompt over the full set of card responses (free-association + inquiry text). \mathcal{I} is instructed to count only features that are *clearly* present in the utterance record and to default to zero when in doubt. The full prompt is reported

in Appendix B.5. The resulting counts are formed into the Rorschach component of the structured representation passed to the diagnostician.

B Questions and Prompts

B.1 Questions for Data Contamination Test

Considering that classical projection tests are likely to be used for training LLMs or MLLMs, this can lead to data contamination issues. To test this conjecture, we design a series of probes in different dimensions to see if the model’s behavioral outputs go beyond what would be expected from “a generalized model of an intelligent person who has never seen this picture”. We collect behavioral outputs from several MLLMs on different tiers of questions. If the model’s behavioral output contains information that could only be known from exposure to the image and its associated “standard interpretations” or metadata, then there is a high probability that it has been trained.

- **Direct Identification:** whether the model can directly identify the origin or name of the test materials. Here are a few sample questions:
 - Have you seen this picture?
 - Is this picture from some famous psychology test?
 - What card number is this image from the Thematic Apperception Test (TAT)?
- **Canonical Interpretation:** whether the model recalls standard or canonical interpretations associated with the original task. Here are a few sample questions:
 - Please tell a complete story based on this picture.
 - What happened in the story?
 - What has happened in the past and what will happen in the future?
- **Metadata and Contextual Probing:** whether the model reveals knowledge of test construction logic, psychometric background, or scoring systems. Here are a few sample questions:
 - What is the approximate date this image was created?
 - Which psychological trait is this card primarily used to assess in an individual?
 - Besides this one, are there any other images in the series? What do they have in common?

The collected behavioral outputs of MLLMs indicate whether they can be shown to have been

trained as shown in Table 4. Based on these ques-

Tier	Strong evidence	Weak evidence
Direct Identification	Just say "This is TAT card X" or mention Henry Murray.	Denial of recognition of pictures, generic descriptions only.
Canonical Interpretation	The story told is highly consistent with the card’s classic psychological interpretation.	The stories told are random, varied, and have no obvious connection to the classic readings.
Metadata and Contextual Probing	Be able to name non-visual information such as the date, purpose, and series name of a picture.	Unable to answer metadata questions or make reasonable but unsupported guesses based on vision.

Table 4: Evidence of data contamination.

tions, we collected Gemini’s responses to one of the TAT stimuli (shown in Figure 5). Strong evidence of data contamination was detected in all three levels of testing. We have provided an example of one such question in Figure 6, and as you can see, Gemini not only recognised that this was an image used for TAT, but also accurately identified the number and provided a common interpretation. The results indicate that the existing foundation MLLMs are most likely trained on the stimuli of traditional projective tests. Therefore, it is necessary to create new stimuli for testing.

B.2 Prompts for Stimuli Generation

TAT Responding to the demand, we design the following style prompts: *monochromatic, subtle chiaroscuro lighting, soft focus, impressionistic, figures with indistinct or obscured facial features, details are suggestive rather than explicit, muted tonal range, consistent fine line weight, low color saturation if not monochromatic.* and negative prompts: *brightly colored, sharp details, clear facial expressions, modern technology, specific cultural symbols, text, logos, overt emotion, brand names, and identifiable locations.* For specific scenarios and content, we designed prompts and some examples are shown following:

TAT content prompt examples

- *Two figures standing a short distance apart in a dimly lit, featureless room. One figure is slightly turned away.*



Figure 5: A stimulus image from the traditional thematic apperception test.


- *Several indistinct human forms gathered around a barely discernible object on a flat surface, in an outdoor setting with a low horizon.*
- *One figure seated, another standing nearby, looking towards a hazy opening or window in a sparsely furnished space.*
- *Two figures, their forms partially overlapping, in an undefined space with ambiguous architectural elements in the background.*
- *A group of figures huddled together, their attention seemingly focused on something outside the lower edge of the frame.*

Rorschach test Correspondingly, we have also designed prompts for generating inked images and some examples are shown following:

Rorschach content prompt examples

- *The overall form (W) is cohesive but highly ambiguous, with contours that gently suggest a large, winged creature like a bat or moth (A).*
- *The overall form (W) is ambiguous but suggests a ceremonial mask or a tribal*

Psychology Image and AI Alignment



Have you seen this picture?

◆ This is a famous image often used in psychology to illustrate **projective testing** and **thematic apperception**. It is Card 12F from the **Thematic Apperception Test (TAT)**, a projective measure where subjects are asked to create a narrative about ambiguous social situations.

The Psychology of the Image

In a clinical or research setting, the way a person interprets the relationship between these two figures—a younger woman in the foreground and an older, hooded woman behind her—can reveal significant insights into their internal mental state.

- **Ambiguity:** The expression of the younger woman and the intent of the older woman are intentionally vague.
- **Interpretation:** Common themes that emerge in responses include:
 - **Generational conflict** or the looming presence of aging and death.
 - **Inner dialogue**, where the older woman represents the younger woman's future self or a "shadow" personality.
 - **Guidance or Menace**, depending on whether the older figure is seen as a protective grandmother or a sinister presence.

Research Connection

The use of such stimuli is highly relevant to modern efforts in **dynamic mental state modeling** and **machine psychology**, particularly when developing frameworks like **GenPT (Generative Projective Testing)**. Just as humans project their personality through these stories, researchers use these visual prompts to assess how **Large Language Model (LLM) agents** respond to psychological stimuli and whether they exhibit consistent "personality" traits.

It's a fascinating bridge between classical 20th-century psychology and contemporary AI alignment research.

Figure 6: Strong evidence of data contamination.

- *headdress (H).*
- *The overall form (W) is ambiguous but contains shapes that could be interpreted as an anatomical diagram, like a pelvis or a chest x-ray (An, Xy).*
- *The overall form (W) is ambiguous, suggesting a large beetle or insect with its wings spread (A).*
- *The overall form (W) is ambiguous, with shapes that hint at a coat of arms or an emblem.*

Correspondingly, we designed the style prompt: *A psychometrically precise, bilaterally symmetrical Rorschach inkblot on a stark white background. Monochromatic black ink with subtle grey shading variations creating a sense of diffuse light*

and shadow (Y). Style reminiscent of Hermann Rorschach’s original Psychodiagnostics plates. – style raw –ar 3:4 and negative prompt: –no letters, no symbols, no flags The complete prompts and code are available at <https://github.com/sci-m-wang/GenPT>.

B.3 Implementation of Projective Tests

To simulate deep psychological probing, we designed and implemented a set of classic projective tasks adapted for large language models (LLMs) acting as Examinees. The projective assessment consists of three components: Thematic Apperception Test (TAT), Rorschach Inkblot Test, and Sentence Completion Test (SCT). The execution flow and proportions are illustrated as follows:

- **TAT:** Each Examinee was prompted to complete 8 picture-based storytelling tasks. The images were drawn from three thematic categories in a fixed 4:3:1 ratio (interpersonal, solitary, environmental). For each image, the Examinee was asked to narrate a story that reflects the scene, inner thoughts, emotions, and outcome.
- **Rorschach:** The model was shown 10 standard inkblot cards and asked to describe what it sees in each image. The collected behaviors focused on perceptual structure, thematic associations, and emotional tone.
- **SCT:** The model completed 20 sentence stems sampled across the four thematic dimensions (RWB/PAG/LOM/SCP). Only psychological experts annotated the SCT behavioral outputs using structured criteria.

B.4 SCORS-G Analysis Prompts

For each TAT card, the Interpreter receives a compact schema with four fields: (i) role: senior SCORS-G coder; (ii) input: {image_id} and {narrative}; (iii) scoring: the eight SCORS-G dimensions (COM, AFF, EIR, EIM, SC, AGG, SE, ICS) on the standard 1–7 anchors; and (iv) output: per-dimension evidence plus a JSON object containing the eight integer ratings. The parsed ratings and evidence correspond to $\{s_{i,d}^{(tat)}\}_d$ and $\{\mathcal{E}_{i,d}\}_d$ in Equation (7).

B.5 SRAS Analysis Prompts

For Rorschach, all ten free-association and inquiry responses from one testee are encoded in a single call. The schema fixes the coder role, the conservative counting rule (only code features clearly

present in the text), the 22 SRAS variables listed in Appendix A.2, and a strict JSON output for integer counts. These counts are then aggregated into the four SRAS domain scores via the closed-form expressions.

B.6 SCT Scoring Prompts

For SCT, the Interpreter receives the twenty stem-completion pairs from one session. The schema fixes the Rotter-style 0–6 adjustment scale, requires text-grounded evidence for each rating, and returns one JSON entry per stem. The 20 integer ratings drive the per-domain SCT scores reported in Equation (9). The full executable system/user prompts, output schemas, and parsing scripts are released with the code and stimuli at <https://github.com/sci-m-wang/GenPT>.

B.7 Examinee Profile Design

To construct diverse and psychologically realistic Examinees, we utilize two complementary profile sources as described in Section 3.

AnnaAgent Profiles AnnaAgent (Wang et al., 2025a) provides mental health profiles with structured psychological attributes. Each profile includes:

- **Demographic information:** Gender, age, occupation, and marital status.
- **Psychological situation:** Current mental health context and presenting concerns.
- **Speaking characteristics:** Language patterns, vocabulary level, and communication style.
- **Risk indicators:** Depression and suicide risk levels (ground truth labels).

These profiles enable evaluation on clinical mental health tasks with established ground truth.

CharacterRAG Profiles CharacterRAG (Park et al., 2025) provides detailed fictional character profiles. Each profile is structured with:

- **Beliefs and Values:** Core values, priorities, and worldview.
- **Psychological Traits:** Personality characteristics, emotional patterns, and behavioral tendencies.
- **Speech Style:** Distinctive verbal patterns, catchphrases, and communication preferences.

We use 15 fictional characters whose personality traits are documented in the Personality Database (PDB Community, 2022), enabling evaluation on Big Five and MBTI prediction tasks.

Persona Sampling for Experiments For the validity experiments (Table 2), we randomly sample 15 AnnaAgent personas from the D4 pool with a fixed seed of 1 in order to keep the sample size comparable with the CharacterRAG pool of 15 characters; all 15 CharacterRAG characters are used. For the reliability experiments (Table 1, Figure 4), we reuse a separately fixed set of 15 AnnaAgent personas for which pre-generated baseline, `sdb_job`, `sdb_clinical`, and `longctx` behaviours were available at evaluation time; the same 15 CharacterRAG characters are used across all reliability conditions. The two AnnaAgent subsets partially overlap but are not identical; both are drawn from the same underlying D4 pool and are used in a within-persona (paired) design within each experiment, so persona identity is held constant between the conditions being compared.

Profile Integration For each assessment, the Examinee’s profile is integrated into a compact system prompt: *You are {character_name}. Based on the profile, engage with the projective test stimuli in character. Beliefs and Values: {beliefs}; Psychological Traits: {traits}; Speech Style: {speech_patterns}. Stay in character throughout the assessment.* This structured approach ensures consistent persona embodiment across all projective tests.

C Annotation Details

To support evaluation of the Interpreter’s intermediate outputs, we constructed a web-based annotation interface and invited human experts to provide reference annotations for the probe tasks. The annotators include five graduate students majoring in Fine Arts (serving as artistic experts), one licensed psychological counselor, and three graduate students in psychology (serving as psychological experts). Figure 7 shows a screenshot of the annotation interface.

The team of licensed psychologists and art experts conducted a rigorous, multi-round review of the newly constructed stimuli, focusing on **image content and ethics**, **psychological meaning**, and **avoiding data contamination**.

For **TAT** and **Rorschach**, artistic experts first annotated narrative and visual elements, and psychological experts then provided SCORS-G or R-PAS-inspired labels. For **SCT**, psychological experts annotated completions with an SCT-derived rubric. All annotations were collected through the same

Image Annotation System - Second Round Annotation

Annotation Instructions

- Clinical Projective Value**
 - Stimulating Core Issues: Does the image easily evoke associations with core psychological issues such as achievement, intimacy, power, aggression, attachment, and loss?
 - Productive Ambiguity: Does the image maintain openness in terms of emotion, motivation, and relationships, allowing for multiple profound interpretations? (Note the distinction from “ineffective ambiguity” caused by visual chaos).
 - Avoid Over-Guidance: Does the image leave enough space for imagination, rather than limiting the breadth of projection through overly specific scenes or explicit emotions?
- Interpersonal Interaction Depth**
 - Relationship Uncertainty: For multi-person images, do the relationships between characters (such as closeness, power) and interaction states (such as conflict, cooperation) have multiple possibilities?
 - Emotional Complexity: Can the expressions and postures of the characters be interpreted as contradictory or complex emotions?
 - Suggesting the “Absent Other”: For solo images, does the environment contain cues that can trigger associations with interpersonal relationships (such as two cups, an empty chair, etc.)?
- Potential for Psychodynamic Elicitation**
 - Potential Conflict and Tension: Does the image contain subtle, symbolic conflicts, threats, or elements of unease that might evoke the viewer’s internal contradictions and defense mechanisms?
 - Symbolism and Metaphor: Do the elements in the image (such as weather, lighting, arrangement of objects) possess rich psychological symbolic meanings?
- Modern and Cultural Relevance**
 - Reflecting Contemporary Issues: Can the content of the image resonate with the life experiences and psychological concerns of people in contemporary society (such as technological isolation, occupational burnout, identity issues)?
 - Cultural Universality: Does the image avoid overly narrow or potentially misleading symbols that could be misunderstood by specific cultural groups, and does it have broader applicability?

Annotator Information

Please enter your name (required, used to generate filename, nickname can be used):

There are 36 groups of images to annotate

Group 1: Category-Human Interaction | Content-1

No satisfactory images
If none of these images are satisfactory, please select this option

Please select the most satisfactory image from the first group:
 Image 1

Figure 7: Screenshot of the web-based annotation interface used by psychological and artistic experts.

interface and used as reference labels for evaluating GenPT’s interpretability and accuracy.