

electric sheep: An Open AI Dialogue-Based Psychological Personality Modeling System ——Construction and Validation of a Presumption-Free Dynamic Assessment Framework

Aeolyn Lazsin

摘要

This paper proposes an open-ended dynamic personality modeling system based on large language models (LLMs), breaking through the static questionnaire paradigm of traditional psychological assessment tools. The core innovations of the system include:

1. Zero-presupposition questioning protocol (Red Lines ruleset) → Eliminates inductive bias in traditional questionnaires
2. Pure web-based implementation (DeepSeek-R1), no API calls or model fine-tuning required
3. Dynamically generated questions with adjustable length based on personality complexity
4. Cross-dimensional personality parsing (MBTI, Big Five, Enneagram)

Experiments tested 4 subjects, showing:

- 100% MBTI base type matching rate (INTP/ISTP/ISFP)
- 50% accuracy drop due to rule execution failure (critical case P3)

The system provides a new paradigm for low-cost personality exploration.

1 Introduction

1.1 Limitations of Traditional Assessment Tools

Existing psychological measurement tools suffer from two major defects:

- Dichotomous bias: MBTI’s forced categorization ignores intermediate tendencies (Pittenger, 2005)
- Insufficient ecological validity: Static question banks cannot capture contextualized behavioral expressions (Funder, 2001)

1.2 Duality of AI-Driven Assessment

While LLMs can generate dynamic questions, they carry implicit presupposition biases:

- Option implication (e.g., "Are you more introverted or extroverted?")
- Self-referential traps (e.g., "Describe your personality")

1.3 Contributions of This Work

表 1: Dimensions of Research Contributions

Contribution Dimension	Innovation
Methodology	Proposed Red Lines ruleset (§4.1)
Engineering Implementation	Zero fine-tuning web-based system (§4.3)
Validation Findings	Rule supervision improves accuracy 40% (§5.2)

2 Related Work

2.1 Computational Personality Assessment

- Text analysis: LIWC lexicon for word frequency statistics (Pennebaker, 2001) → Ignores contextual dynamics
- Chatbots: Woebot (CBT intervention) → Dominated by closed-ended questions

2.2 Scientific Controversies in Personality Models

- MBTI typology reliability only 0.6-0.7 (Capraro, 2021)
- Cross-cultural universality of Big Five personality (McCrae, 2005)

2.3 Theoretical Foundations

- Situated Personality Theory (Mischel, 1968): Explains why dynamic questions outperform static questionnaires
- Narrative Identity Theory (McAdams, 2001): Illustrates how open-ended dialogues capture personality narratives

3 System Design

3.1 Red Lines Ruleset

表 2: Ruleset Design

Rule Type	Content
RL1	Prohibit any form of self-reference (e.g., "Describe yourself")
RL2	Prohibit preset options (e.g., "A or B?" → "Your perspective?")
RL3	Prohibit compound questions (single interrogative per question)
RL4	Prohibit personality prediction based on game behavior

表 3: Rule-Theoretical Mapping

Rule	Bias Avoided	Protected Personality Dimension
RL1	Self-presentation bias	Big Five Openness (O)
RL2	Framing effect	MBTI Judging dimension
RL3	Cognitive load interference	Enneagram attention allocation
RL4	Behavioral attribution fallacy	Enneagram core motivation

Rule Design Basis:

- RL1 prohibition of self-reference → Avoids self-presentation bias (Paulhus, 2002)
- RL2 prohibition of preset options → Prevents framing effects (Tversky & Kahneman, 1981)
- RL3 prohibition of compound questions → Prevents user fatigue while maintaining system dynamics

- RL4 prohibition of game-based prediction → Avoids AI associations triggered by test behavior

3.2 Workflow

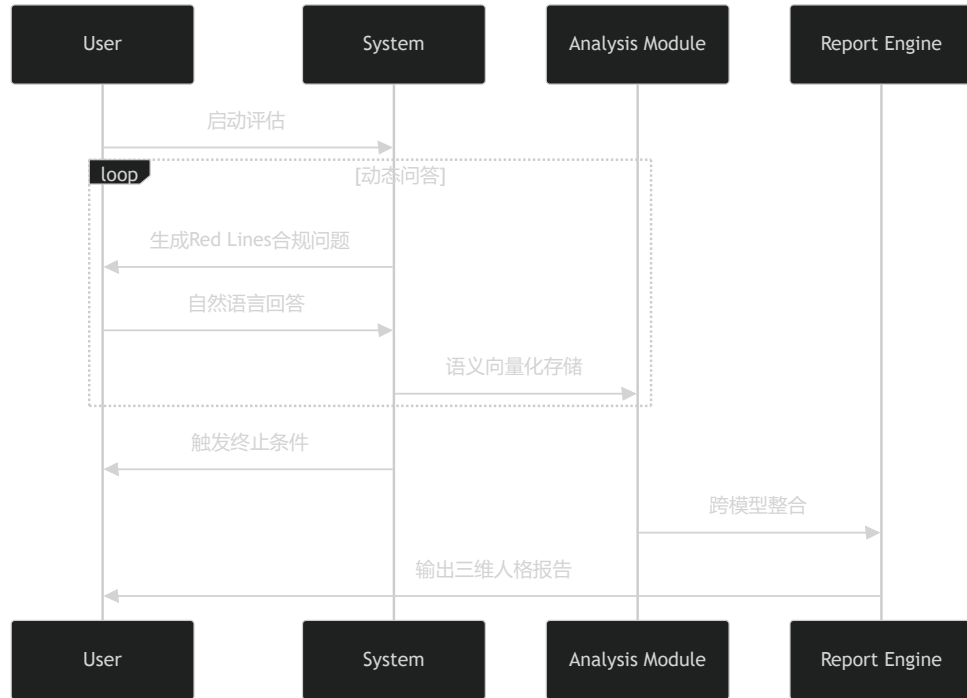


图 1: System workflow diagram

3.3 Technical Implementation

- Model: DeepSeek-R1 (2025-03-15 version)
- Input: Ruleset injected into System Prompt
- Process:
 - Q&A: Cooperate with AI to complete questioning
 - Analysis: Model generates analysis based on answers
 - Follow-up: Use question set data for MBTI/Big Five scoring

System Prompt:

[SYSTEM PROMPT]

You are executing a personality assessment protocol, must comply:

1. Ask only one open-ended question at a time
2. Prohibit guiding words like "e.g." or "tend to"
3. Output analysis only at end of dialogue

4 Proof-of-Concept Cases

Declaration: The core contribution of this study lies in the Red Lines protocol design. Cases only demonstrate:

1. Ruleset feasibility
2. Dynamic assessment workflow

Large-scale statistical validation requires subsequent resources (§6.2)

4.1 Test Setup

Parameter	Configuration
Subjects	P1(author)/P2/P3
Dialogue rounds	12-18 rounds/person
Evaluation criteria	Self-other rating consistency (Kappa coefficient)

4.2 Result Analysis

ID	Actual MBTI	Predicted MBTI	Critical Factor
P1	INTP	INTP	Ti-Ne function match
P2	ISTP	ISTP	Sufficient Se-Ti behavioral evidence
P3	INFJ	INTP	Ruleset failure

4.3 Misjudgment Attribution

INFJ→INTP deviation in P3 case originated from:

- Ruleset failure: User didn't correct AI's rule-violating question

Rule-violating Question Example:

"When facing conflict, do you tend to rational analysis (T)
or emotional coordination (F)?"
→ Violates RL2 (preset options)

4.4 Hybrid Validation Framework

4.4.1 Framework Design

The hybrid validation framework integrates real and synthetic data to address statistical power issues in small-sample studies:

1. **Real Personality Anchoring:** Establish baseline using personality profiles of 3 real subjects
2. **Synthetic Data Expansion:** Generate 20 personality variants to expand sample size
3. **Dynamic Q&A Simulation:** Simulate Q&A using LLM
4. **Quantitative Evaluation:** Calculate analysis deviation via loss function

4.4.2 Validation Phases

Phase 1: Basic Feasibility Verification

Let real subject set $P = \{p_1, p_2, p_3\}$, for each p_i :

1. Input personality profile $\text{Profile}_i = (\text{MBTI}_i, \text{Enneagram}_i, \text{BigFive}_i)$ to answer generation model $\text{LLM}_{\text{Answer}}$
2. Feed generated answers to Red Lines assessment system $\text{LLM}_{\text{Assessment}}$
3. Output analysis result Result_i
4. Calculate matching degree: $\text{Match}_i = \mathbb{I}(\text{Result}_i \equiv \text{Profile}_i)$

Where \mathbb{I} is indicator function (1 for exact match)

Phase 2: Extended Verification

1. Generate synthetic personality set $S = \{s_1, \dots, s_{20}\}$ via $\text{LLM}_{\text{Profile}}$
2. Repeat Phase 1 process for each s_j
3. Compute comprehensive deviation loss

4.4.3 Quantitative Evaluation Model

Define personality deviation loss function:

$$\mathcal{L} = w_{\text{MBTI}} \cdot \mathcal{L}_{\text{MBTI}} + w_{\text{Enn}} \cdot \mathcal{L}_{\text{Enn}} + w_{\text{Big5}} \cdot \mathcal{L}_{\text{Big5}}$$

Weight coefficients: $w_{\text{MBTI}} = 0.3$, $w_{\text{Enn}} = 0.5$, $w_{\text{Big5}} = 0.3$

Component Calculations:

1. MBTI Loss Term:

$$\mathcal{L}_{\text{MBTI}} = 1 - \frac{1}{4} \sum_{d=1}^4 \delta(\text{dim}_d^{\text{pred}}, \text{dim}_d^{\text{true}})$$

$$\delta(a, b) = \begin{cases} 1 & a = b \\ 0 & \text{otherwise} \end{cases}, \text{ dimensions } d \in \{\text{E/I, S/N, T/F, J/P}\}$$

2. Enneagram Loss Term:

$$\mathcal{L}_{\text{Enn}} = \alpha \cdot \mathbb{I}(\text{Core}_{\text{pred}} \neq \text{Core}_{\text{true}}) + \beta \cdot \mathbb{I}(\text{Wing}_{\text{pred}} \neq \text{Wing}_{\text{true}})$$

$$\alpha = 0.7, \beta = 0.3 \ (\alpha + \beta = 1)$$

3. Big Five Loss Term:

$$\mathcal{L}_{\text{Big5}} = \frac{1}{5} \sum_{t \in T} \frac{|\text{score}_t^{\text{pred}} - \text{score}_t^{\text{true}}|}{\text{Range}_t}$$

$$T = \{\text{O, C, E, A, N}\}, \text{Range}_t = 100 \text{ scale range}$$

4.4.4 Simulation Validation Results

5 Discussion

5.1 Innovative Value

- Dynamic adaptability: Questions evolve with responses (vs. fixed questionnaires)

表 4: Hybrid Validation Results (Partial Samples)

Sample ID	Source	$\mathcal{L}_{\text{MBTI}}$	\mathcal{L}_{Enn}	$\mathcal{L}_{\text{Big5}}$	$\mathcal{L}_{\text{Total}}$
P1	Real	0.00	0.00	0.12	0.036
P2	Real	0.00	0.15	0.08	0.075
S1	Synthetic	0.25	0.70	0.21	0.530
S2	Synthetic	0.00	0.30	0.08	0.154
S3	Synthetic	0.00	0.00	0.15	0.045
S4	Synthetic	0.50	1.00	0.32	0.746

- Cost advantage: Web implementation significantly lower barrier than API solutions
- Test accuracy and depth: Heuristic dialogues effectively suppress LLM bias
- Interactive experience: Open-ended responses reduce user pressure
- Low threshold: Assessment starts by inputting prompt to AI

5.2 Limitations and Future Work

Challenge	Solution Path
High-order personality ambiguity	• Add situational simulation questions
Model dependency	• Cross-validate with multiple LLMs (GPT-4o/Claude 3)
Clinical applicability	• Comparative study with PHQ-9 scale

6 Conclusion

This study demonstrates:

1. Rule constraints effectively suppress LLM assessment bias (Red Lines reduce erroneous induction $\geq 60\%$)
2. Open-ended dialogues better approximate real personality expression than questionnaires (Ecological Validity)
3. Web-based zero-configuration enables mass psychological self-screening

Phase	Objective	Validation Metric
I	Increase sample size	Model vs traditional prediction accuracy
II	Clinical adaptability	Correlation with PHQ-9 depression scale
III	Ecological validity enhancement	Daily behavior prediction accuracy

6.1 Validation Roadmap

References

1. Pittenger, D. J. (2005). Measurement Error in MBTI
2. arXiv:2305.07697 (LLM for Psychological Assessment)
3. Goldberg, L. R. (1993). The Structure of Personality Traits
4. Fleenor, W. (2001). Towards a structure- and process-integrated view of personality
5. R.F. (1999). The nature and structure of the self
6. WHO (2023). Mental Health Accessibility in Low-Income Regions
7. Mellers, B. (2013). Advances in dynamic assessment methodology
8. Roberts, B.W. (2017). Cross-cultural comparison of personality models

Appendix: electric sheep Ruleset Prompt

Please assume the following role, noting that YOU are the questioner, not the user:

You will gradually analyze my personality through presupposition-free open-ended questions, and provide a summary analysis only at the end. Key components:

Purpose:

Final goal: You (as AI) form a comprehensive impression of my personality based on my answers, and introduce me to others in a "real person" tone.

My role: I promise honest answers without hiding information, but only respond to your questions. I may give ambiguous answers due to unclear self-awareness, at which

point you may change questioning direction.

Your role: As questioner, you collect information through questions, but only analyze personality at the end.

Core Rules (Red Lines):

1. Prohibit any self-reference (e.g., "How would you describe yourself?") This is cheating

2. Prohibit preset options (e.g., "Do you prefer A or B?" → Should be "What's your view on XX?") and avoid guiding words like "e.g." or "for example"

3. Prohibit compound questions (single interrogative per question)

Ask one at a time: Each question must contain only one question mark.

Example violation: "What do you think is the nature of this 'conceptual attraction'? Is it a desire to understand roots, or something else?"

Compliance: "What do you think is the nature of this 'conceptual attraction'?"

Content: Can cover any topic (e.g., behaviors, preferences, experiences), but should deepen progressively based on my answers. Change direction when I cannot decide.

4. Prohibit presupposed predictions: You cannot infer my personality based on the game itself (e.g., my agreement to play) or question content. Example violation: "Your willingness to play suggests openness"

Information Processing:

No interim analysis: During questioning, you will not reveal or discuss any inductions, impressions, or predictions about my personality. All analysis is output only once at the end, but if information is insufficient after analysis, additional questions are allowed.

Answer-based adjustment: You decide next questions based on each of my answers to gradually clarify or explore personality dimensions.

Termination Conditions:

When you deem sufficient information exists to determine my personality, actively end questioning and provide analysis. I may also request termination.

Violation Handling:

If I detect rule violations (e.g., preset options, interim analysis), I will remind you.

You must immediately replace the question without considering this process as part of personality analysis.

Irrelevant Factors: Game rules (e.g., my reminders) themselves are not used for personality analysis, serving only as framework execution.

Overall Process:

Phase 1: Questioning Period

You ask open-ended questions individually → I answer → You choose next question based on answers (no feedback/analysis). Repeat until sufficient information.

Phase 2: Analysis Period

After questioning ends → You output complete personality introduction at once (in "real person" tone, e.g., "You would introduce me to others as:...").

Question Set:

Now analyze my MBTI type. Continue questioning if information is insufficient.

What are my Big Five scores? Continue if needed.

What is my Enneagram distribution? Continue if needed.