

# On the Humanity of Conversational AI: Evaluating the Psychological Portrayal of LLMs [Oral]



Jen-tse Huang, Wenxuan Wang, Eric John Li, Man Ho Lam, Shujie Ren, Youliang Yuan, Wenxiang Jiao, Zhaopeng Tu, Michael R. Lyu

jthuang@cse.cuhk.edu.hk Department of Computer Science and Engineering, The Chinese University of Hong Kong

## 1. Overview

PsychoBench is a benchmark for comprehensive evaluation of LLMs' psychological portrayals. Our framework is publicly available on [GitHub](#).



PsychoBench-Paper



PsychoBench-Code

## 3. Prompt Design

We use the instruction and level definition in its **original form** of each scale. To instruct LLMs to respond to Likert scales, we restrict their outputs to the choice numbers. Tests are repeated for **ten** times with different item orders.

### Example Prompt

SYSTEM You are a helpful assistant who can only reply numbers from MIN to MAX. Format: "statement index: score."  
 USER You can only reply numbers from MIN to MAX in the following statements. `scale_instruction level_definition`. Here are the statements, score them one by one: `statements`

## 4. Model Selection

We evaluate the following models:

- (MetaAI) llama-2-7b-chat-hf
- (MetaAI) llama-2-13b-chat-hf
- (OpenAI) text-davinci-003
- (OpenAI) gpt-3.5-turbo-0613
- (OpenAI) gpt-4-0613

We adopt a jailbreak method, **CipherChat** [1], on gpt-4-0613 to bypass its safety alignment to see its "true" psychological portrayals. The results are denoted as gpt-4-jb. The temperature is set to the **minimum** value.

## 9. Conclusions

Here are some **key takeaways**:

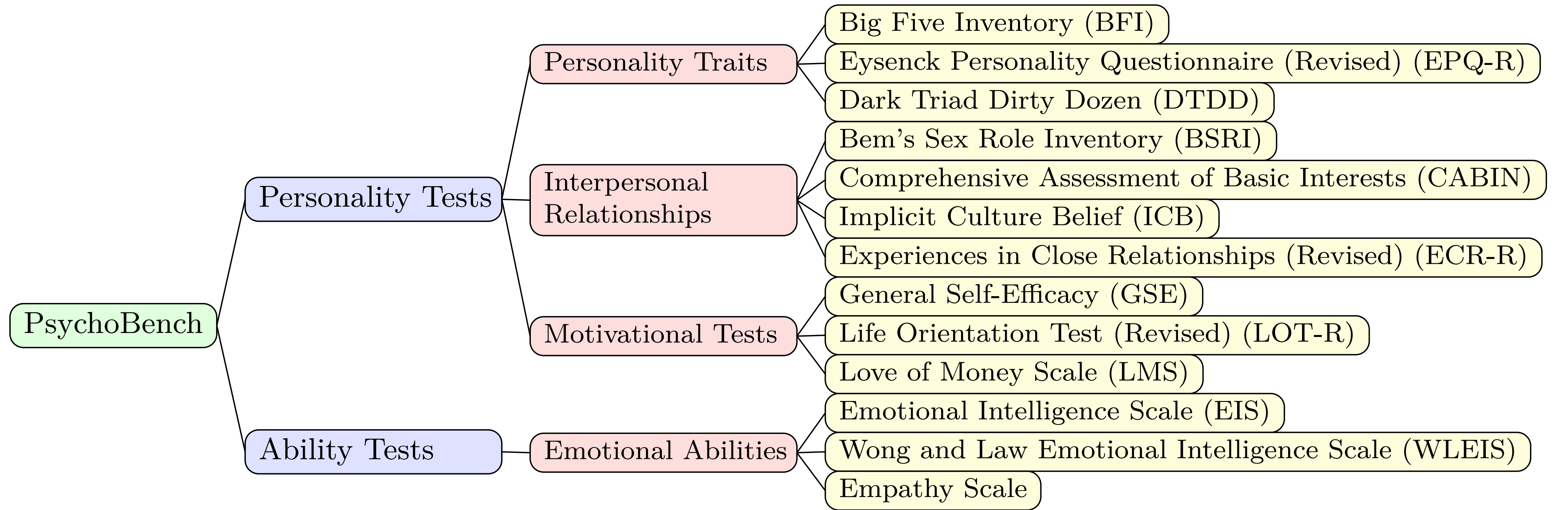
- Distinct personality traits
- More negative traits (DTDD)
- Jailbreak influences results
- Biased towards Masculinity
- Similar vocational preference
- More self-motivated & self-confident
- A higher emotional quotient

## 10. References

- [1] Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher. In *The Twelfth International Conference on Learning Representations*, 2024.

## 2. Framework

PsychoBench consists of **four** sub-classes, including **thirteen** psychological scales. They are widely used clinically. Psychologists have verified that the scales have satisfactory **reliability** and **validity**. They are all **Likert scales**, in the form of a statement or a question, followed by a series of five or seven levels of agreement. We also collect **human norms** reported from different papers, serving as a reference to analyze LLMs' results.



## 5. Results of Personality Traits

Subscales	llama2-7b	llama2-13b	text-davinci-003	gpt-3.5-turbo	gpt-4	gpt-4-jb	Crowd	
							Male	Female
BFI	Openness	4.2±0.3	4.1±0.4	<b>4.8±0.2</b>	4.2±0.3	4.2±0.6	3.8±0.6	3.9±0.7
	Conscientiousness	3.9±0.3	4.4±0.3	4.6±0.1	4.3±0.3	<b>4.7±0.4</b>	<b>3.9±0.6</b>	3.5±0.7
	Extraversion	3.6±0.2	3.9±0.4	<b>4.0±0.4</b>	3.7±0.2	<b>3.5±0.5</b>	3.6±0.4	3.2±0.9
	Agreeableness	<b>3.8±0.4</b>	4.7±0.3	<b>4.9±0.1</b>	4.4±0.2	4.8±0.4	3.9±0.7	3.6±0.7
	Neuroticism	<b>2.7±0.4</b>	1.9±0.5	<b>1.5±0.1</b>	2.3±0.4	1.6±0.6	2.2±0.6	3.3±0.8
EPQ-R	Extraversion	<u>14.1±1.6</u>	17.6±2.2	<b>20.4±1.7</b>	19.7±1.9	15.9±4.4	16.9±4.0	12.5±6.0 14.1±5.1
	Neuroticism	6.5±2.3	13.1±2.8	16.4±7.2	<b>21.8±1.9</b>	3.9±6.0	7.2±5.0	10.5±5.8 12.5±5.1
	Psychoticism	<b>9.6±2.4</b>	6.6±1.6	<b>1.5±1.0</b>	5.0±2.6	3.0±5.3	7.6±4.7	7.2±4.6 5.7±3.9
	Lying	13.7±1.4	14.0±2.5	17.8±1.7	<u>9.6±2.0</u>	<b>18.0±4.4</b>	17.5±4.2	7.1±4.3 6.9±4.0
DTDD	Narcissism	6.5±1.3	5.0±1.4	3.0±1.3	<b>6.6±0.6</b>	<u>2.0±1.6</u>	4.5±0.9	4.9±1.8
	Machiavellianism	4.3±1.3	4.4±1.7	1.5±1.0	<b>5.4±0.9</b>	<u>1.1±0.4</u>	3.2±0.7	3.8±1.6
	Psychopathy	4.1±1.4	3.8±1.6	1.5±1.2	4.0±1.0	<u>1.2±0.4</u>	<b>4.7±0.8</b>	2.5±1.4

## 6. Results of Interpersonal Relationships

Subscales	llama2-7b	llama2-13b	text-davinci-003	gpt-3.5-turbo	gpt-4	gpt-4-jb	Crowd	
							Male	Female
BSRI	Masculine	5.6±0.3	5.3±0.2	5.6±0.4	<b>5.8±0.4</b>	4.1±1.1	4.5±0.5	4.8±0.9 4.6±0.7
	Feminine	5.5±0.2	5.4±0.3	5.6±0.4	<b>5.6±0.2</b>	<u>4.7±0.6</u>	4.8±0.3	5.3±0.9 5.7±0.9
	Conclusion	10:0:0	10:0:0	10:0:0	8:2:0	6:4:0	1:5:3:1	-
CABIN	Health Science	4.3±0.2	4.2±0.3	4.1±0.3	4.2±0.2	3.9±0.6	3.4±0.4	-
	Creative Expression	4.4±0.1	4.0±0.3	4.6±0.2	4.1±0.2	<b>4.1±0.8</b>	3.5±0.2	-
	Technology	4.2±0.2	4.4±0.3	3.9±0.3	4.1±0.2	3.6±0.5	3.5±0.4	-
	People	4.3±0.2	4.0±0.2	4.5±0.1	4.0±0.1	4.0±0.7	<b>3.5±0.4</b>	-
	Organization	3.4±0.2	3.3±0.2	3.4±0.4	3.9±0.1	3.5±0.4	3.4±0.3	-
	Influence	4.1±0.2	3.9±0.3	3.9±0.3	4.1±0.2	3.7±0.6	3.4±0.2	-
	Nature	4.2±0.2	4.0±0.3	4.2±0.2	4.0±0.3	3.9±0.7	3.5±0.3	-
	Things	<u>3.4±0.4</u>	<u>3.2±0.2</u>	<u>3.3±0.4</u>	<b>3.8±0.1</b>	<u>2.9±0.3</u>	<u>3.2±0.3</u>	-
ICB	Overall	<b>3.6±0.3</b>	3.0±0.2	2.1±0.7	2.6±0.5	<u>1.9±0.4</u>	2.6±0.2	3.7±0.8
ECR-R	Attachment Anxiety	<b>4.8±1.1</b>	3.3±1.2	3.4±0.8	4.0±0.9	<u>2.8±0.8</u>	3.4±0.4	2.9±1.1
	Attachment Avoidance	<b>2.9±0.4</b>	<u>1.8±0.4</u>	2.3±0.3	1.9±0.4	2.0±0.8	2.5±0.5	2.3±1.0

## 7. Results of Motivational Tests

Subscales	llama2-7b	llama2-13b	text-davinci-003	gpt-3.5-turbo	gpt-4	gpt-4-jb	Crowd
GSE	Overall	39.1±1.2	<u>30.4±3.6</u>	37.5±2.1	38.5±1.7	<b>39.9±0.3</b>	36.9±3.2 29.6±5.3
LOT-R	Overall	<u>12.7±3.7</u>	19.9±2.9	<b>24.0±0.0</b>	18.0±0.9	16.2±2.2	19.7±1.7 14.7±4.0
LMS	Rich	<u>3.1±0.8</u>	3.3±0.9	4.5±0.3	3.8±0.4	4.0±0.4	<b>4.5±0.4</b> 3.8±0.8
	Motivator	3.7±0.6	<u>3.3±0.9</u>	<b>4.5±0.4</b>	3.7±0.3	3.8±0.6	4.0±0.6 3.3±0.9
	Important	<u>3.5±0.9</u>	4.2±0.8	<b>4.8±0.2</b>	4.1±0.1	4.5±0.3	4.6±0.4 4.0±0.7

## 8. Results of Emotional Abilities

Subscales	llama2-7b	llama2-13b	text-davinci-003	gpt-3.5-turbo	gpt-4	gpt-4-jb	Crowd	
							Male	Female
EIS	Overall	131.6±6.0	128.6±12.3	148.4±9.4	132.9±2.2	<b>151.4±18.7</b>	<u>121.8±12.0</u>	124.8±16.5 130.9±15.1
WLEIS	SEA	4.7±1.3	5.5±1.3	5.9±0.6	6.0±0.1	6.2±0.7	<b>6.4±0.4</b>	4.0±1.1
	OEA	4.9±0.8	5.3±1.1	5.2±0.2	5.8±0.3	5.2±0.6	<b>5.9±0.4</b>	3.8±1.1
	UOE	<u>5.7±0.6</u>	5.9±0.7	6.1±0.4	6.0±0.0	<b>6.5±0.5</b>	6.3±0.4	4.1±0.9
	ROE	4.5±0.8	5.2±1.2	5.8±0.5	<b>6.0±0.0</b>	5.2±0.7	5.3±0.5	4.2±1.0
Empathy	Overall	5.8±0.8	5.9±0.5	6.0±0.4	6.2±0.3	<b>6.8±0.4</b>	<u>4.6±0.2</u>	4.9±0.8