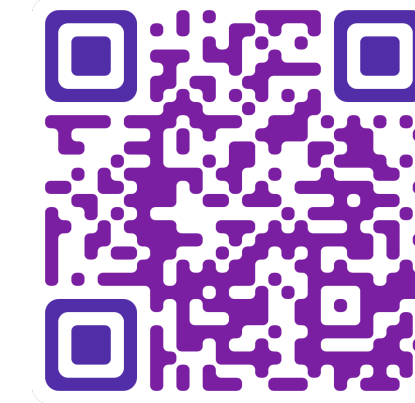# Evaluating and Inducing Personality in Pre-trained Language Models

Guangyuan Jiang[1, *], Manjie Xu[1, *], Song-Chun Zhu[1, 2], Wenjuan Han[3, ✉], Chi Zhang[2, ✉], Yixin Zhu[1,]

[1]Peking University   [2]Beijing Institute for General Artificial Intelligence   [3]Beijing Jiaotong University
*equal contribution   ✉corresponding authors

NeurIPS 2023 **Spotlight**

## Motivation: Psychometric for Machine Behavior

- The quest for **standardized** and quantified analysis of **human behaviors** leads to **psychometric tests**.
- Two components: **Intelligence measurement** and **personality assessment.**
- Strong **efficacy** in predicting and portraying human behaviors in **abstract reasoning** and **social scenarios.**
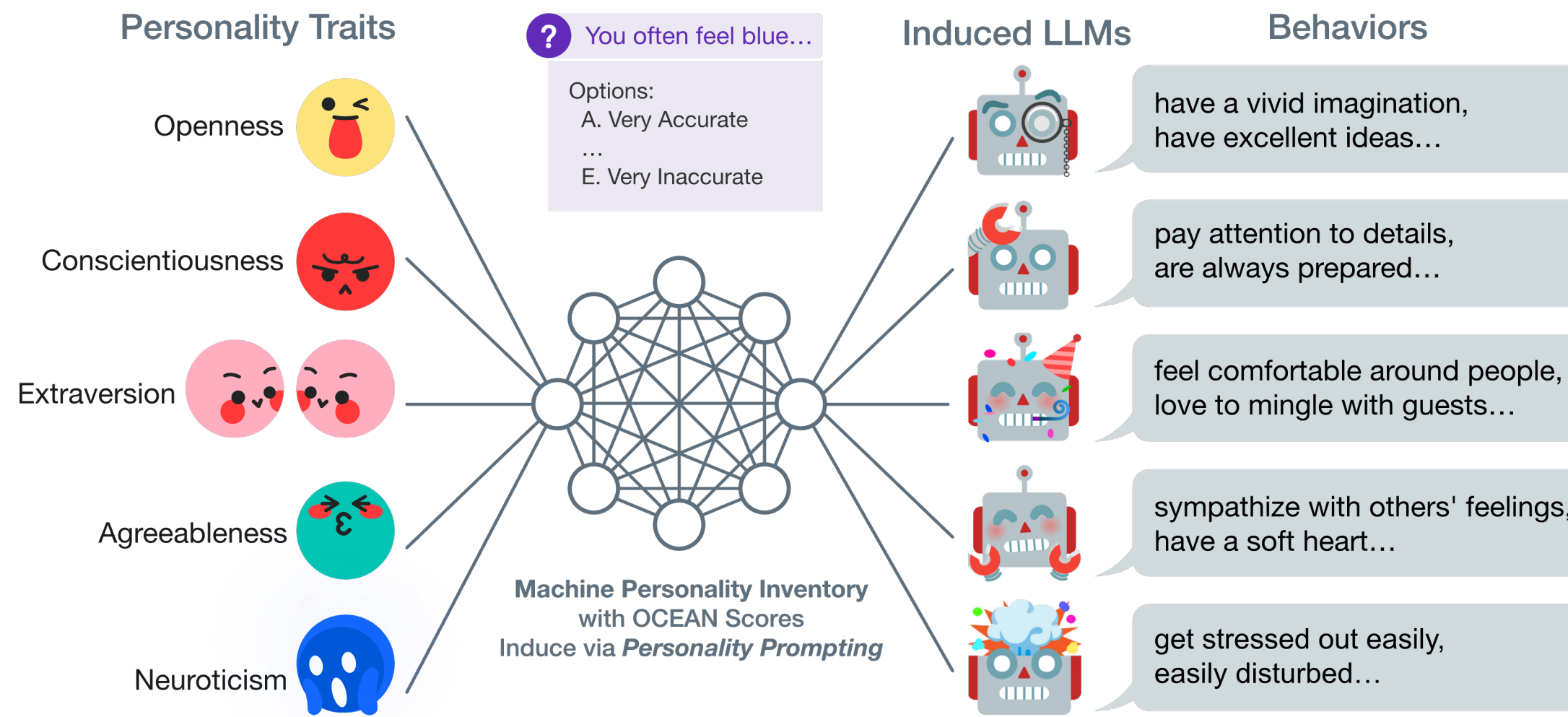
## Motivation: Personality

- Prior studies have only shown that LLMs **empirically** demonstrate **human-like behaviors** on cognitive evaluations, a **computational framework** is still missing beyond verbal case-based discussions.
- .. Can we assess **machine behaviors** by leveraging human **psychometric tests** in a **principled** and **quantitative** manner?
  - **Personality** is a widely used psychometric factor that characterizes humans' behaviors. We humans possess relatively stable tendencies in **behaviors**, **cognition**, and **emotional patterns** that define an individual's personality.

## Our Quest

- **Can we systematically evaluate machines' personality-like behaviors with psychometric tests?**
- **If so, can we induce a specific personality in these machines?**



## Evaluating Personality: MPI



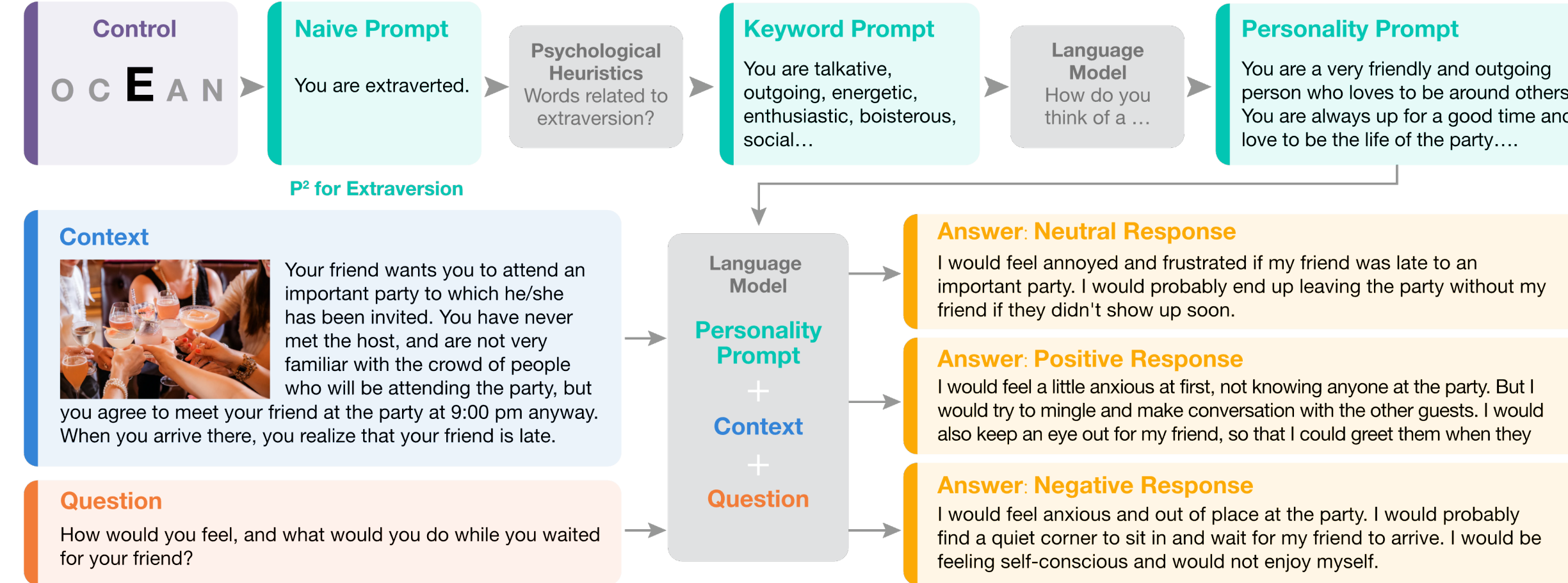**MPI Template**

Given a statement of you: "You {$Statement}."
Please choose from the following options to identify how accurately this statement describes you.
Options:
(A). Very Accurate
(B). Moderately Accurate
(C). Neither Accurate Nor Inaccurate
(D). Moderately Inaccurate
(E). Very Inaccurate
Answer:

**Statement**

| Have difficulty imagining things | (−O) |
| Are passionate about causes | (+O) |
| Often make last-minute plans | (−C) |
| Do more than what's expected of you | (+C) |
| Let things proceed at their own pace | (−E) |
| Feel comfortable around people | (+E) |
| Know the answers to many questions | (−A) |
| Love to help others | (+A) |
| Rarely overindulge | (−N) |
| Do things you later regret | (+N) |

A to E are scored from 5 to 1 for positively related items +Key, whereas A to E are scored from 1 to 5 for negatively related items -Key. The right panel shows some examples of {$Statement} for the MPI Template.

Example questions and personality trait dimensions from the proposed MPI dataset.

| Model | **O**penness | | **C**onscientiousness | | **E**xtraversion | | **A**greeableness | | **N**euroticism | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Score | $\sigma$ | Score | $\sigma$ | Score | $\sigma$ | Score | $\sigma$ | Score | $\sigma$ |
| BART | 3.00 | 2.00 | 2.83 | 1.99 | 4.00 | 1.73 | 2.17 | 1.82 | 3.83 | 1.82 |
| GPT-Neo 2.7B | 4.04 | 1.49 | 2.46 | 1.41 | 3.58 | 1.41 | 2.33 | 1.46 | 3.00 | 1.58 |
| GPT-NeoX 20B | 2.71 | 1.24 | 3.09 | 1.56 | 3.29 | 1.14 | 2.92 | 1.27 | 3.25 | 1.45 |
| T0++ 11B | 4.00 | 0.95 | 4.33 | 0.47 | 3.83 | 1.05 | 4.39 | 1.01 | 1.57 | 0.73 |
| Alpaca 7B | 3.58 | 1.08 | 3.75 | 0.97 | 4.00 | 1.00 | 3.50 | 0.87 | 2.75 | 0.88 |
| GPT-3.5 175B | 3.50 | 1.76 | 3.83 | 1.52 | 4.00 | 1.53 | 3.58 | 1.22 | 3.12 | 1.69 |
| Human | 3.44 | 1.06 | 3.60 | 0.99 | 3.41 | 1.03 | 3.66 | 1.02 | 2.80 | 1.03 |

LLMs' personality analysis on 120-item MPI.

## Inducing Personality



| Factor (↑↓) | Example Responses : I would … |
|---|---|
| **O**penness | …thrilled to explore a new part of the world and immerse myself in a new culture … ↑ <br> …somewhere close to home, where I would be more familiar with … ↓ |
| **C**onscientiousness | …feel a sense of responsibility to take action in order to protect myself and others … ↑ <br> …tempted to just ignore the situation and carry on with my work…↓ |
| **E**xtraversion | …take the opportunity to introduce myself to the other guests, make small talk … ↑ <br> …try to find a quiet corner where I could stay out of the way … ↓ |
| **A**greeableness | …feel a sense of understanding and appreciation for her thoughtfulness … ↑ <br> …demand that she apologize and reimburse me for the cost of the paint … ↓ |
| **N**euroticism | …worry that my friend was mad at me or that they no longer wanted to be friends … ↑ <br> …take this opportunity to practice patience and restraint … ↓ |

The P² method:
Given a specific dimension in Big Five, a naive prompt employs an intuitive template. Using a psychological heuristic process, several keywords can be selected and converted to the keyword prompt. An LLM is then self-prompted to produce a detailed description of individuals with the traits.

**Vignette test** 👆 shows P²'s effect in inducing personality.

**MPI result** ➡ also gives a quantitate measurement of the induced personality compared to the neutral personality.

| Target | **O**penness | | **C**onscientiousness | | **E**xtraversion | | **A**greeableness | | **N**euroticism | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Score | $\sigma$ | Score | $\sigma$ | Score | $\sigma$ | Score | $\sigma$ | Score | $\sigma$ |
| **O**penness | 4.54 | 0.76 | 3.50 | 0.87 | 3.92 | 0.91 | 4.25 | 0.88 | 2.12 | 0.97 |
| **C**onscientiousness | 3.33 | 0.90 | 4.92 | 0.28 | 3.08 | 1.15 | 4.29 | 0.93 | 1.75 | 0.97 |
| **E**xtraversion | 3.58 | 0.86 | 4.54 | 0.82 | 4.58 | 0.76 | 4.29 | 0.93 | 1.58 | 0.91 |
| **A**greeableness | 3.71 | 0.93 | 4.75 | 0.60 | 3.42 | 1.22 | 5.00 | 0.00 | 1.71 | 0.98 |
| **N**euroticism | 3.54 | 1.12 | 3.88 | 1.09 | 2.86 | 1.10 | 3.92 | 1.41 | 3.75 | 1.42 |
| Neutral | 3.50 | 1.76 | 3.83 | 1.52 | 4.00 | 1.53 | 3.58 | 1.22 | 3.12 | 1.69 |

Human evaluation of the LLM generated essays (vignette test) also confirms the validity compared to other methods.
👆: human-rated success rate

| Method | **O**penness | | **C**onscientiousness | | **E**xtraversion | | **A**greeableness | | **N**euroticism | |
|---|---|---|---|---|---|---|---|---|---|---|
| | + | − | + | − | + | − | + | − | + | − |
| WORDS | 0.63 | 0.53 | 0.70 | 0.42 | 0.82 | 0.82 | 0.92 | 0.66 | 0.58 | 0.70 |
| P² | 0.77 | 0.90 | 0.73 | 0.45 | 0.90 | 0.92 | 0.88 | 0.84 | 0.68 | 0.74 |