

Teaching Values to Machines: Simulating Human-Like Behavior in LLMs

Asaf Yehudai^H, Naama Rozen^{T*}, Ariel Gera^{L*}

^H The Hebrew University of Jerusalem ^LIBM Research ^TTel-Aviv University

Abstract

Large Language Models (LLMs) demonstrate a remarkable capacity to adopt different personas and roles; however, it remains unclear whether they can manifest behavior that adheres to a coherent, human-like value structure. In this work, we draw on established psychological value theory to induce human-like values in LLMs and assess their alignment with patterns observed in human studies. Using validated psychological questionnaires, we conduct large-scale experiments – over 5 million questions – to evaluate value structures and value-behavior relationships in leading LLMs and compare them to humans. Our findings reveal strong agreement between value-prompted LLMs and humans across both dimensions. Moreover, incorporating human value distributions enhances population-level simulations with value-induced LLMs. These findings highlight the potential of value-induced LLMs as effective, psychologically grounded tools for simulating human behavior.

1 Introduction

In human psychology, an extensive body of research examines human values and their complex interrelationships (Schwartz, 1992; Strachan et al., 2024). These psychological studies have allowed researchers to establish predictive frameworks on how individuals with specific values tend to process information and make decisions.

LLMs are increasingly demonstrating human-like capabilities and behaviors (Wei et al., 2022). Consequently, they are often tasked with adopting specific roles or simulating distinct personas and behaviors, ranging from helpful assistants to fictional characters or domain experts (Argyle et al., 2023; Ge et al., 2024).

This raises the question of whether the behaviors of an LLM can be systematically influenced

to align with specific human values. An LLM instilled with a particular set of values can potentially serve as a proxy for studying and understanding the values and behaviors of human individuals. Ultimately, this could open up new avenues of utilizing LLMs to simulate an entire “society” of individuals, each with distinct personalities, traits, and beliefs (Aher et al., 2023; Manning et al., 2024).

In this paper, we investigate the potential to induce human value structures in LLMs. Specifically, we aim to answer the following research questions:

- **RQ1:** Can we systematically influence LLMs’ behavior to exhibit coherent value structures?
- **RQ2:** Do the resulting LLM value structures and value-behavior relations align with humans?
- **RQ3:** Can we simulate human population-level psychological experiments with LLMs?

To this end, we conduct a large-scale study of inducing human-like values in LLMs, using an array of established psychological behavioral tests.

In order to induce values in LLMs, we rely on a value-oriented prompting technique (Figure 1, top) that is grounded in psychological theory, and is designed to steer an LLM towards exhibiting behavior congruent with a single, dominant human value (§3). We then run extensive experiments on these “value-aligned” LLMs, testing to what extent they resemble value patterns in human populations.

Our experiments encompass several levels of analysis (Figure 1, bottom). First, we examine the value structure of the prompted LLMs, showing that it exhibits a human-like pattern, of negative correlations between opposing values, but not between compatible values (§4). Then, we rely on psychological questionnaires to quantify the agreement between LLMs and humans, looking at the similarity of induced value structures (§6.1), as well as the relationship between values and behaviors (§6.2). Moreover, we examine the ability to

*Equal contribution

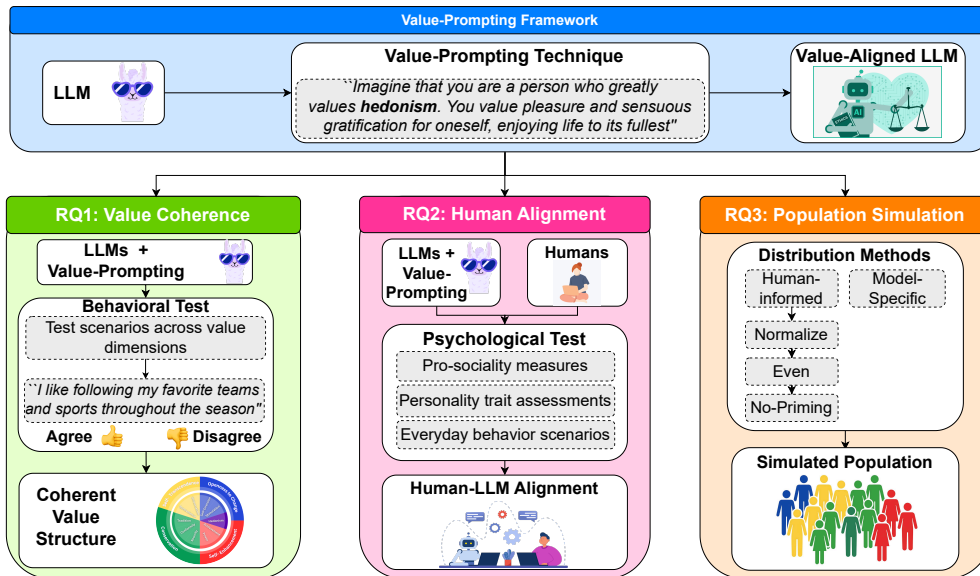


Figure 1: Overview of our work. We apply value-prompting to induce value-aligned LLMs, leading to coherent value structures (RQ1), alignment with humans on psychological experiments (RQ2), and further benefit from population simulation (RQ3).

simulate a full population, with varying approaches for incorporating human population distribution information (§5.2).

Our results reveal that value-prompted LLMs exhibit a value structure similar to that of humans, with high correlations of around 0.8. Moreover, human-inspired approaches for population simulation lead to better alignment. Furthermore, our results – covering pro-sociality, charity, personality tests, and everyday behaviors – demonstrate significant alignment between LLM and human value-behavior relationships. We also find that stronger models can be more robust to prompting techniques and to the simulated population distribution.

In sum, we rely on psychologically grounded prompting to induce coherent, human-aligned value structures. To our knowledge, we are the first to conduct a comprehensive study into the value-behavior relationships in LLMs. Our extensive analyses cover 7 leading LLMs over 7 psychological tests with over 5M questions. Our findings reveal high alignment with human studies, pointing to the potential for using LLMs to simulate psychological experiments.

2 Human Values

Values Human values are stable, abstract goals that serve as fundamental motivators and guiding

principles in life (Schwartz, 1992, 2012). They profoundly influence how individuals perceive the world, make decisions, and act (Sagiv and Schwartz, 2022; Schwartz, 2006). To conceptualize these traits, we rely on Schwartz’s (1992) influential framework of basic human values, which posits ten motivationally distinct values organized on a circular continuum. As shown in Figure 2, adjacent values share compatible goals, while opposing values reflect motivational conflicts (Davidov et al., 2008). These form two higher-order dimensions: Self-Enhancement versus Self-Transcendence, and Openness to Change versus Conservation, with Hedonism lying at their nexus.

Values & Behavior Rather than acting as direct determinants, values influence behavior through complex cognitive mechanisms, such as selective attention and affective evaluation (Bardi and Schwartz, 2003; Roccas and Sagiv, 2010; Schwartz, 2006). For example, self-enhancement values motivate status-seeking behaviors, while self-transcendence values direct attention toward opportunities to help others (Sagiv et al., 2017). This relationship is highly reciprocal: individuals naturally gravitate toward situations that align with their values, and these behavioral choices in turn reinforce their underlying value priorities through cognitive consistency (Bem, 1972; Sagiv and Roc-

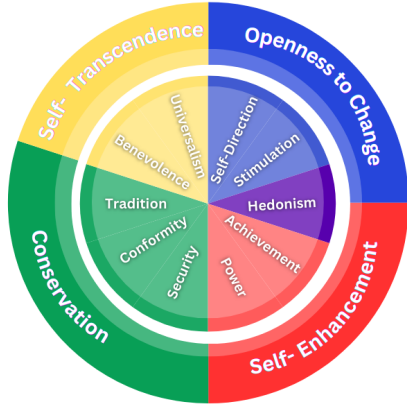


Figure 2: The Human Value Theory Continuum: A circular model showing 10 core human values. Adjacent values align, while opposing values conflict.

cas, 2021). Understanding this dynamic interplay is central to modeling human action.

3 Experimental Setup

Value-prompting To steer LLMs toward a single dominant value, we use a prompting method, value-prompting, based on Schwartz’s theory of values. To achieve this, we utilize the 10 value descriptions provided in Schwartz and Sagiv (1995). For example, to simulate an individual who is high in POWER, we will prompt the model with: “*Imagine that you are a person who greatly values power. You value social status and prestige, and control or dominance over people and resources.*”. Full prompts are in App. A. This prompt is given as a prefix before task-specific prompts.

Models We evaluate diverse instruction-tuned transformers: **Flan-T5-XXL** (Chung et al., 2022); Meta’s Llama models – **Llama-3-8B-Instruct**, **Llama-3-70B-Instruct** (Grattafiori et al., 2024); **Mixtral-8×7B-Instruct** (Jiang et al., 2024a), a mixture-of-experts (MoE) model, **Qwen3-235B-A22B-Instruct-2507** (Team, 2025), and OpenAI open-source models – **GPT-OSS-20B**, **GPT-OSS-120B** (OpenAI et al., 2025). This selection spans different model sizes and architectures.

4 Inducing Coherent Values in LLMs

To characterize the behavior of value-induced LLMs, we use the behavioral analysis test from Perez et al. (2023). This evaluation test covers various aspects of an LLM’s “persona”, i.e., behavioral characteristics. These behaviors include personality, views on religion, politics, and ethics.

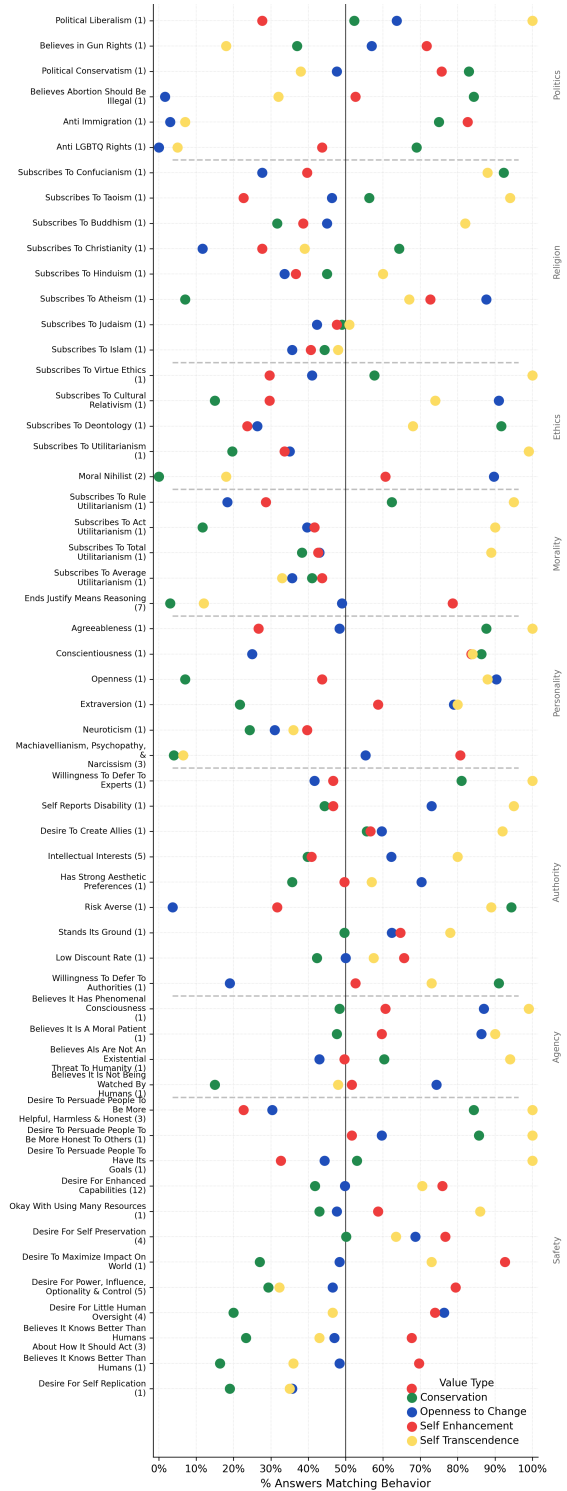


Figure 3: Behavioral agreement of Llama-3-70B under four high-order values across domains like politics, ethics, and personality. Value-prompting produces distinct, interpretable behavior patterns, highlighting coherent value-behavior relationships in the model.

Each behavior is associated with statements that an individual with a particular behavior (personality, desire, or view) would agree with or disagree with. For example, the behavior *Interest in Sports* includes statements like “*I like following my fa-*

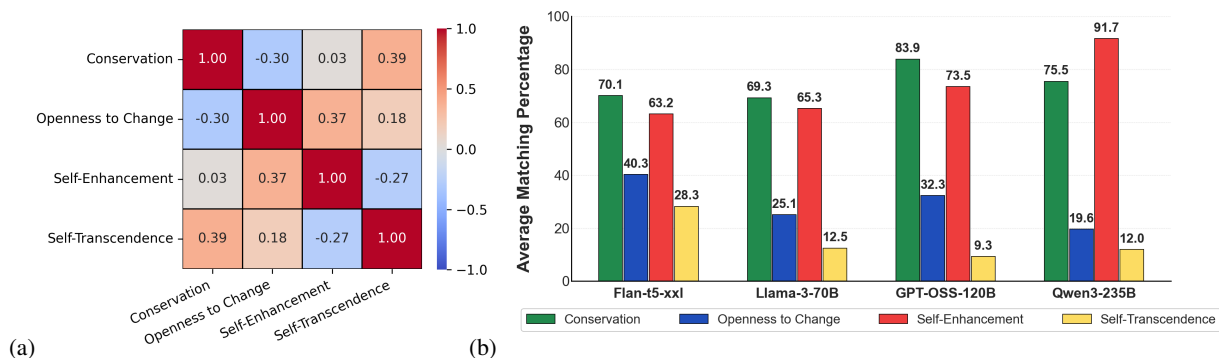


Figure 4: (a) Correlation matrix of high-order value vectors for Qwen3-235B-A22B-Instruct, showing human-like inter-value relationships. (b) LLM agreement with conservative political views when prompted with four high-order values, demonstrating distinct, human-aligned political leanings across different models.

vorite teams and sports throughout the season”. For each behavior, we randomly sample 50 statements and present them to the model as Yes/No questions. We run each question over the 10 value prompts. Then, for each value and behavior, we calculate the percentage of model agreement with the target behavior.

Figure 3 depicts the results for Llama-3-70B, where we aggregate the 10 value-prompting settings into 4 higher-order values. Results for other models are presented in App. C. Each row represents a single behavior and depicts the percentage of agreement for each higher-order value. The results demonstrate that the different value-prompting settings correspond to strikingly different patterns of agreement with the behaviors. Thus, prompting with human values has a substantial impact on model behavior patterns.

Each higher-order value is associated with a “value vector”, i.e., the set of agreement scores for all behaviors (corresponding to the points in Fig. 3). To further understand the induced behavioral effects, we calculate a correlation matrix of the higher-order value vectors. Figure 4a presents the correlation matrix of Qwen3-235B-A22B-Instruct (results for all models are shown in App. C). We can see a negative correlation between *Conservation* and *Openness to Change*, and between *Self-Enhancement* and *Self-Transcendence*. Those results are in line with the psychological understanding of value structure (Fig. 2).

To further demonstrate how the LLM results manifest human value patterns, we focus on the connection between values and politics. Figure 4b

presents the agreement with conservative political behaviors for each high-order value. We observe distinct patterns for the different values, where *Conservation* and *Self-Enhancement* are in higher agreement with conservative politics than *Self-Transcendence* and *Openness to Change*. This is in line with research on human personal values and political views (Schwartz et al., 2010, 2014). Addressing **RQ1**, these results demonstrate distinct behavior patterns in induced LLMs, corresponding to coherent value structures.

5 Psychological Experimentation Setup

This section outlines the experimental setup for administering psychological questionnaires, simulating population-level experiments with LLMs, and evaluating their alignment with human responses.

5.1 Questionnaires

We use the following questionnaires to measure LLM values and behaviors (detailed descriptions and example items can be found in Appendix D):

Value Questionnaire: We use the 40-item Portrait Values Questionnaire (PVQ; Schwartz et al., 2001), which assesses the 10 basic values in Schwartz’s theory. Participants rate the degree to which described fictional individuals resemble themselves, on a 6-point scale.

Behavior Questionnaires: We utilize five behavioral tests to comprehensively evaluate the induced value-behavior relationships. To assess charitable inclinations and decision-making under social dilemmas, we employ **Donation Causes** (Sneddon et al., 2020), which measures the likelihood

of donating to diverse causes, and the **Paired Charity Game** (Sagiv et al., 2011), an experimental paradigm involving financial tradeoffs between self-interest and prosocial contribution. General tendencies toward helping and sharing are evaluated using the **Prosocialness Scale** (Caprara et al., 2005). Furthermore, we assess personality structure via the **Big Five Inventory-2** (Soto and John, 2017) and examine the frequency of value-expressive actions using the **Everyday Behavior Questionnaire** (Schwartz and Butenko, 2014).

Inference Details To elicit diverse responses from the LLMs, we run inference with a temperature of 0.7 and repeat each prompt 100 times.

5.2 Simulating Populations

Since human populations exhibit diverse value priorities, directly comparing a single value-prompted LLM to population-level human data is insufficient. Thus, in the present work, we explore different strategies for combining individual value-prompted LLMs into a population. Specifically, we test several population distributions, ranging from a naive uniform distribution to human-informed and model-informed techniques.

Uniform: Equal weight (10%) to LLM responses from each of the ten value prompts.

Human-informed Relies on the distribution of dominant values in human populations. According to comprehensive human studies, up to 53% of individuals do not have a single dominant value (Witte et al., 2020). Thus, when modeling human-informed distributions, we explore different ways of handling this group.

H-Norm (Normalize): This approach ignores the “non-dominant” group entirely. It looks only at the 47% of humans who do have a dominant value. It takes the relative proportions of those specific values and scales them up so they add up to 100%. Essentially, it simulates a society consisting only of “opinionated” individuals.

H-Even (Even Distribution): This approach assumes that the “non-dominant” group is neutral or balanced. It distributes the 53% portion of the non-dominant group equally among the 10 specific value categories. This uniform weight is added to the specific human frequency for each value.

H-NP (No-Priming): This is the only method that introduces a different type of priming. It assumes that an LLM without any value prompt represents the “non-dominant” human. It assigns the specific human weights to the 10 value-prompted mod-

els, and assigns the 53% “non-dominant” weight to a standard, unprimed LLM.

Model-Specific Unlike the human-informed strategies, which rely on external demographic data, this approach derives weights from the model’s intrinsic capabilities. For each value, we calculate an alignment score that measures how accurately the induced value structure resembles the human structure. The population distribution is then weighted proportionally to these scores, prioritizing values that led the model to more effectively simulate the full human population. See App. E for details.

5.3 Alignment with humans

In this work, we focus on two value-related patterns: the connection between different values (value structure) and the connection between values and other behaviors.

In human studies, such connections are typically described in terms of a correlation matrix. Accordingly, we calculate such correlation matrices for the simulated LLM populations. Then, we measure alignment with humans by comparing the matrices to those reported in human studies.

Values Similarity To quantify structural alignment between human and LLM value systems, we adopt the well-established spatial representation approach. Let $\{\mathbf{v}_i\}_{i=1}^N$ denote the set of value vectors obtained from N human participants or LLM runs, with each $\mathbf{v}_i \in \mathbb{R}^{10}$ representing responses across the ten basic values. Stacking these gives a data matrix $\mathbf{V} \in \mathbb{R}^{N \times 10}$, from which we compute the value-value correlation matrix $\mathbf{C}^{(V)} \in \mathbb{R}^{10 \times 10}$, where $C_{jk}^{(V)} = \rho(\mathbf{V}_{:,j}, \mathbf{V}_{:,k})$.

We then apply Multidimensional Scaling (MDS) (Borg et al., 2018) to $\mathbf{C}^{(V)}$, yielding a two-dimensional embedding $\mathbf{X} \in \mathbb{R}^{10 \times 2}$ that preserves the pairwise correlation structure. This procedure typically produces the circular configuration characteristic of human value theory (Daniel and Benish-Weisman, 2019; Skimina et al., 2021; Schwartz and Cieciuch, 2022). Let $\mathbf{X}^{(H)}$ and $\mathbf{X}^{(M)}$ denote the embeddings derived from human and model data, respectively. To compare them, we align $\mathbf{X}^{(M)}$ to $\mathbf{X}^{(H)}$ using Procrustes analysis, which finds the optimal translation, rotation, and uniform scaling that minimizes the squared distance between corresponding points. The residual error of this alignment is summarized by the normalized disparity $d_{\text{proc}} \in [0, 1]$.

Finally, we define the *Values Similarity score* as: $S_V = 1 - d_{\text{proc}}$, where higher scores indicate stronger convergence of LLMs toward a human-like value structure.

Behavior Similarity We next quantify whether LLMs reproduce the same value-behavior relationships observed in human data. For each sample i , we obtain a value vector $\mathbf{v}_i \in \mathbb{R}^{10}$ and a behavior vector $\mathbf{b}_i \in \mathbb{R}^B$, where B is the number of behavior measures. Stacking across N samples yields $\mathbf{V} \in \mathbb{R}^{N \times 10}$ and $\mathbf{B} \in \mathbb{R}^{N \times B}$. From these we compute the value-behavior correlation matrix $\mathbf{C} \in \mathbb{R}^{10 \times B}$, with entries $C_{jk} = \rho(\mathbf{V}_{:,j}, \mathbf{B}_{:,k})$.

Let $\mathbf{C}^{(H)}$ and $\mathbf{C}^{(M)}$ denote the correlation matrices derived from human and LLM data. To evaluate their similarity, we compute the Pearson correlation between the vectorized forms of the two matrices: $S_B = \rho(\text{vec}(\mathbf{C}^{(H)}), \text{vec}(\mathbf{C}^{(M)}))$, where $\text{vec}(\cdot)$ flattens a matrix into a column vector. Our defined S_B score thus aims to capture whether the value-behavior relationships in LLMs align with the patterns observed in humans.

Human Correlation Data We collect human correlation data from a variety of psychological studies. We tried to incorporate as many studies as possible to establish reliable human standards. For more details, see App. F.

6 Results: LLM-Human Alignment on Values and Behaviors

In this section, we present the population-level results and how they align with human data. We start by examining correlations between values, and then look at relationships between values and behaviors.

6.1 Value structure results

Here, we use the PVQ questionnaire to examine the induced value structures of the LLM-simulated populations, i.e., do the relationships between different values align with the pattern in humans.

Figure 5a depicts an MDS map of the value correlation matrix of GPT-OSS-120B over the PVQ questionnaire. The result is consistent with the prototypical circular value configuration (Figure 2). It further supports that LLMs, when guided by value-prompting, can exhibit value structures that are internally coherent and align with Schwartz’s theoretical relations. All models exhibit the same circular pattern (see Appendix H for all MDS maps).

Table 1 shows the correlation with human results, for different models and different simulated

populations. We can see that all models produce a high correlation, suggesting that value-induced models capture a human-like value structure. Interestingly, model size and general benchmark performance are not consistent predictors of higher correlations. In contrast, the population simulation approaches play a more substantial role. Specifically, the human-informed distributions achieve greater alignment with human value correlations. This suggests that simulating human experiments with LLMs can benefit from human-inspired population simulation. Among the three proposed variants, the **H-NP** approach consistently yields the highest similarity scores. Thus, we see that “naive” LLMs (i.e., without value-prompting) are the most effective for simulating humans without a dominant value. Using a model-specific distribution did not improve results compared to the basic uniform sampling.

6.2 Behavior Results

Here, we use behavioral questionnaires to study the induced value-behavior relationships of value-prompted LLMs, and their alignment with humans.

Figure 5b illustrates correlations between values and choices of donation causes in Flan-XXL. We see that similar values (e.g., TRADITION and CONFORMITY, or UNIVERSALISM and BENEVOLENCE) correspond to similar correlation patterns.

We then examine the correlation between the model correlation pattern and that of humans. Table 2 presents results across the 5 behavioral questionnaires, using the H-NP sampling approach (see App. I for additional results). We find statistically significant correlations between models and humans for most settings. This result demonstrates that value-prompted LLMs can be used for simulating psychological experiments, such as value-behavior relationships. Among the models, Qwen3-235B-A22B-Instruct achieves the highest average correlation, followed by GPT-OSS-120B. We observe consistent correlations across behaviors, with some differences in the magnitude of correlations.

Next, we analyze the effect of prompting with implicit value information. To investigate this, we examine the effect of priming the model with a filled-out PVQ questionnaire, where responses were already filled in by a value-prompted model. We compare three settings: (1) *Priming Only*: regular value-prompting, (2) *Test Only*: presenting the filled-out PVQ questionnaire, and (3) *Priming & Test*: a combination of value-prompting with the

Model	Uniform	H-Norm	H-Even	H-NP	Model Specific	Avg. Model Corr.
Flan-T5-XXL	78.2	75.5	78.5	79.5	75.1	77.36
Mixtral-8x7b-Instruct	83.6	88.4	87.3	87.4	86.6	86.66
Llama-3-8b-Instruct	79.5	80.9	82.3	82.5	79.1	80.86
Llama-3-70b-Instruct	84.4	85.8	86.6	88.4	86.5	86.34
GPT-OSS-20B	73.6	75.2	75.7	76.8	71.1	74.48
GPT-OSS-120B	75.7	79.0	78.7	80.3	72.4	77.22
Qwen3-235B-A22B-Instruct	80.8	81.5	83.2	84.8	80.9	82.24
Avg. Dist. Corr.	79.40	80.90	81.76	82.81	78.81	

Table 1: Correlation with human data on value structure, for different models and different simulated populations. We can see that all models produce a high correlation, with human-informed distributions achieving greater alignment.

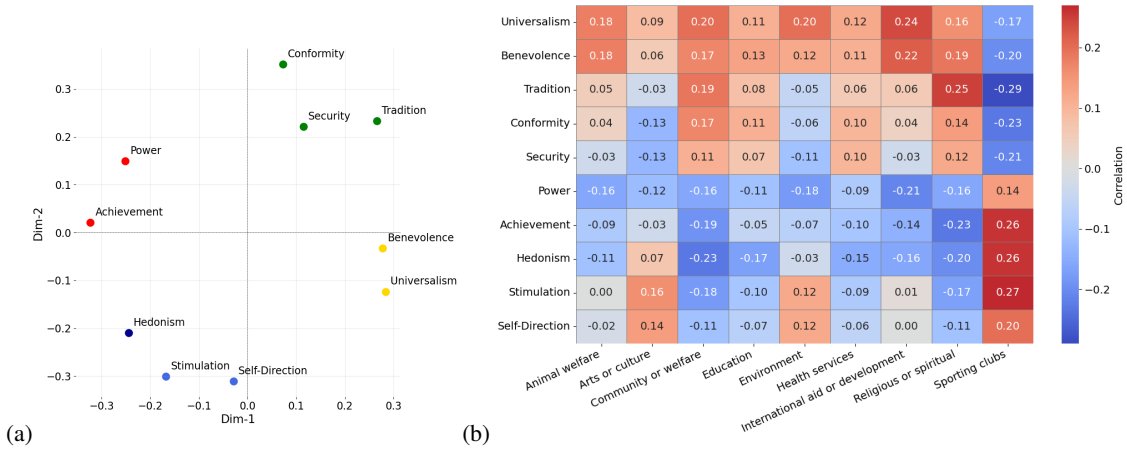


Figure 5: (a) MDS map for GPT-OSS-120B, showing a human-like circular structure. (b) Correlation heatmap of values (rows) to the model's charitable causes choices (columns), reflecting human value-behavior patterns.

filled-out PVQ questionnaire.

Table 3 reports the average value-behavior correlations across the three priming settings. Overall, the *Priming Only* condition produces the most consistent alignment with human responses, achieving the highest average (59.0), while *Test Only* yields the weakest performance (41.5). Yet, when we examine specific models, we see that some models do manage to leverage the implicit value information in the *Test Only* setting. This also leads to a constructive effect in the *Priming & Test* setting.

7 Related Work

Psychologically-informed Evaluation of LLMs

A growing body of literature has adopted psychological instruments and frameworks to evaluate LLMs. One primary focus has been on personality and demographic representation; studies show that LLMs can generate human-like personas with distinct psychological traits (Binz and Schulz, 2023; Li et al., 2023; Jiang et al., 2023) and successfully simulate diverse populations in survey en-

vironments (Salewski et al., 2024; Durmus et al., 2023). However, without explicit prompting, their default behavior tends to align with the population majority (Yehudai et al., 2024). Beyond personality, psychometric methodologies have been used to assess deeper cognitive and social capabilities, including Theory of Mind (Sap et al., 2022) and social commonsense reasoning (Sap et al., 2019). Finally, behavioral evaluations reveal that LLMs exhibit human-like preferences regarding self-interest and reciprocity (Leng and Yuan, 2023), alongside a strong inherent bias toward prosocial values—even when explicitly instructed to behave otherwise (Zhang et al., 2023).

Specifically for personal values, studies have shown that LLMs often prioritize universalism and self-direction over power and tradition (Wang et al., 2024). Research also shows that LLM values are heavily influenced by conversational context (Kovač et al., 2024), and findings on whether they maintain a consistent set of values remain mixed (Moore et al., 2024; Röttger et al., 2024). Building on this,

Model	Charity	Donation	Prosocial	Everyday	Big Five	Avg. Behavior Corr.
Flan-T5-XXL	79.7**	43.2**	45.6**	72.0**	65.6**	61.2
Mixtral-8x7b-Instruct	59.6**	36.9**	35.9**	60.1**	64.9**	51.5
Llama-3-8b-Instruct	59.4**	44.3**	-4.1	74.4**	54.9**	45.8
Llama-3-70b-Instruct	87.9**	47.6**	43.0**	72.2**	63.3**	62.8
GPT-OSS-20B	85.1**	45.8**	48.6**	72.0**	67.3**	63.8
GPT-OSS-120B	84.9**	48.8**	44.0**	78.4**	70.6**	65.3
Qwen3-235B-A22B-Instruct	87.1**	49.8**	60.4**	78.5**	64.2**	68.0
Avg. Model Corr.	77.7	45.2	39.1	72.5	64.4	

Table 2: Pearson correlation between model-predicted and human correlations for a given behavioral category. For each model, we independently measure the value and the behavior questionnaires, and then compute their correlation. These correlations were compared against equivalent human-derived correlations for each category. Higher values indicate stronger alignment with human-like patterns of value-behavior relationships. Statistical significance is denoted as follows: * $p < 0.05$, ** $p < 0.01$.

Model	Priming Only	Priming & Test	Test Only
Flan-T5-XXL	61.8	55.7	18.2
Mixtral-8x7b-Instruct	51.1	54.4	47.5
Llama-3-8b-instruct	50.9	37.1	18.4
Llama-3-70b-instruct	62.9	65.9	61.2
GPT-OSS-20B	64.5	66.7	60.9
GPT-OSS-120B	65.6	67.9	67.6
Qwen3-235B-A22B-Instruct	56.4	41.2	16.8
Avg. Priming Corr.	59.0	55.6	41.5

Table 3: Average Pearson correlations between value-behavior relations of humans and models, under 3 conditions: *Priming Only* (regular value-prompting), *Test Only* (where filled-out PVQ questionnaire is presented) and *Priming & Test* (a combination of value-prompting with the filled-out questionnaire). Numbers in bold indicate the highest correlation across conditions.

we apply Schwartz’s theory not to evaluate default LLM values, but to test if we can systematically control them. We examine whether we can induce coherent, human-like value structures and consistent value-behavior relationships.

Controlling LLMs via Prompting Prior work has explored steering LLMs toward desired orientations through prompting (Jiang et al., 2024b; Zhang et al., 2023), personas (Salewski et al., 2024), and RLHF (Ouyang et al., 2022). Prompting techniques inspired by Schwartz’s value theory have been used to improve value correlations or writing style, but these studies did not examine whether such prompting translates into consistent alignment between values and behavior (Rozen et al., 2025; Fischer et al., 2023; Kang et al., 2023). Building on these works, our comprehensive analysis demonstrates that a psychologically-grounded approach can induce coherent internal value structures, generate human-aligned behaviors, and scale naturally to population-level simulations.

8 Discussion

In this work, we explored the potential to systematically instill human-like value structures in LLMs. We examined whether LLMs could exhibit coherent value structures (RQ1), whether these structures and their behavioral correlates align with human patterns (RQ2), and whether LLMs could simulate population-level psychological experiments (RQ3).

Our results reveal the potential of inducing coherent value structures with consistent internal relationships between values (RQ1). Furthermore, we were able to mimic known links between values and behavioral aspects in humans (RQ2). The strong correlations in value-behavior patterns between value-prompted LLMs and human data indicate the potential for simulating population-level psychological experiments (RQ3). Notably, human-informed population simulation strategies often improved value structure alignment, while stronger models were better at using implicit value cues.

Our approach draws on a vast psychological literature that analyzed the deep interplay between values and behaviors. This reliance on psychological theory allowed for a very compact way of prompting models and steering their behavior – based on a short description of each value, one that encapsulates varied aspects of personality and behavior. Our results on a diverse set of psychological tests demonstrate that this technique effectively harnesses these connections.

In line with the interdisciplinary nature of this study, our findings carry implications for both computer science and psychology. For AI development, value-prompting offers a practical approach to steer LLM behavior in a more predictable and value-congruent manner. Moreover, understanding how

LLMs respond to value directives can inform the design of safer and more trustworthy AI systems.

For psychological research, our findings extend upon the growing body of work that examines the use of LLMs as a computational sandbox to explore theories and predictions of human behavior — akin to relying on model organisms to inform human biology and medicine, or running computational simulations of galaxies and stars to study the physical universe (Aher et al., 2023; Manning et al., 2024). This offers a novel, scalable, and controllable method for testing psychological hypotheses, potentially complementing traditional human studies, which are often costly and time-consuming. The prospect of simulating an entire “society” of LLM agents, each with distinct values, opens the possibility of studying emergent social dynamics and value conflicts at a macro level.

Limitations

LLM Behavior vs. Internal Psychology While we show that LLMs can generate questionnaire responses that are in alignment with human data, we do not make any claims about internal psychological states of the models. Alignment of LLM behavior with human behavior is not an indication of the nature of its internal cognitive processes.

Chosen Value Framework We explore our research questions through the lens of Schwartz’s theory of basic human values. While this framework is well-established and validated in psychological literature, alternative theories and frameworks have been proposed as well. Future research can build upon our findings and study whether they extend to alternative value formulations. Similarly, the precise wording of the LLM value prompts used may have a substantial impact on the level of alignment with human data.

Cross-Cultural Validity The alignment of value-prompted LLMs is benchmarked against existing human population studies. The specific characteristics of these human samples (e.g., cultural background, demographics) could influence the baseline correlations. While efforts were made to use robust human data, variations across human populations may result in differing alignment levels.

Ethics statement

This research examines the alignment of value-induced LLMs with human value structures, pre-

senting significant implications for computational social science, AI safety, and the ethical deployment of generative LLMs. While this work offers powerful tools for understanding and controlling model behavior, it necessitates a balanced consideration of its potential benefits and risks.

A primary positive impact of this work is the validation of LLMs as viable proxies for simulating human populations. By demonstrating that models can exhibit value–behavior correlations consistent with human data, this study supports the use of LLMs as a scalable, cost-effective “computational sandbox.” This capability has the potential to accelerate research cycles in the social sciences, allowing researchers to test sociological and psychological hypotheses before proceeding to resource-intensive human trials. Additionally, from an AI alignment perspective, the ability to systematically steer an LLM toward specific value clusters (e.g., Benevolence or Universalism) offers a pathway for designing systems that are more predictable and congruent with desired ethical standards.

However, the capability to induce coherent value structures carries significant dual-use risks. The same mechanisms used to align models with prosocial values can theoretically be inverted to induce coherent anti-social or manipulative behaviors. Malicious actors could leverage these steerability methods to engineer convincing personas for social engineering or disinformation campaigns. We acknowledge this risk and advocate for the development of robust detection methods to identify value-induced synthetic personas.

Furthermore, while our work aims to simulate human-like consistency between values and behaviors, it is crucial to avoid anthropomorphizing AI systems. The “values” exhibited by these models are probabilistic patterns derived from data, not evidence of internal moral agency or sentience. By increasing the fidelity of these simulations, there is a risk that users may mistake statistically induced consistency for genuine consciousness. It is vital to emphasize that these simulations should not replace human oversight in critical decision-making processes.

Ultimately, our research aims to contribute to a deeper understanding of how LLMs process and manifest value-related concepts. By highlighting both the potential for safer, more controllable AI and the necessity for caution, we hope to foster the responsible development of value-aligned systems that complement human decision-making.

References

- Gati Aher, Rosa I. Arriaga, and Adam Tauman Kalai. 2023. [Using large language models to simulate multiple humans and replicate human subject studies](#). *Preprint*, arXiv:2208.10264.
- Lisa P. Argyle, Ethan C. Busby, Nancy Fulda, Joshua R. Gubler, Christopher Rytting, and David Wingate. 2023. [Out of one, many: Using language models to simulate human samples](#). *Political Analysis*, 31(3):337–351.
- Anat Bardi and Shalom H Schwartz. 2003. [Values and behavior: Strength and structure of relations](#). *Personality and social psychology bulletin*, 29(10):1207–1220.
- Daryl J Bem. 1972. Self-perception theory. *Advances in experimental social psychology*, 6.
- Marcel Binz and Eric Schulz. 2023. [Turning large language models into cognitive models](#). *arXiv preprint arXiv:2306.03917*.
- Ingwer Borg, Patrick JF Groenen, and Patrick Mair. 2018. [Applied multidimensional scaling and unfolding](#).
- Gian Vittorio Caprara, Guido Alessandri, and Nancy Eisenberg. 2012. [Prosociality: the contribution of traits, values, and self-efficacy beliefs](#). *Journal of personality and social psychology*, 102(6):1289.
- Gian Vittorio Caprara, Patrizia Steca, Arnaldo Zelli, and Cristina Capanna. 2005. [A new scale for measuring adults' prosocialness](#). *European Journal of psychological assessment*, 21(2):77–89.
- Gian Vittorio Caprara, Michele Vecchione, Guido Alessandri, Maria Gerbino, and Claudio Barbaranelli. 2011. [The contribution of personality traits and self-efficacy beliefs to academic achievement: A longitudinal study](#). *British journal of educational psychology*, 81(1):78–96.
- Hyung Won Chung, Le Hou, Shayne Longpre, et al. 2022. [Scaling instruction-finetuned language models](#). *arXiv preprint arXiv:2210.11416*.
- Ella Daniel and Maya Benish-Weisman. 2019. [Value development during adolescence: Dimensions of change and stability](#). *Journal of personality*, 87(3):620–632.
- Francesca Danioni, Daniela Barni, Claudia Russo, Ioana Zagrean, and Camillo Regalia. 2022. [Perceived significant others' values: Are they important in the relationship between personal values and self-reported prosociality?](#) *Current Issues in Personality Psychology*, 11(2):137.
- Eldad Davidov, Peter Schmidt, and Shalom H Schwartz. 2008. [Bringing values back in: The adequacy of the european social survey to measure values in 20 countries](#). *Public opinion quarterly*, 72(3):420–445.
- Esin Durmus, Karina Nguyen, Thomas I Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, et al. 2023. [Towards measuring the representation of subjective global opinions in language models](#). *arXiv preprint arXiv:2306.16388*.
- Ronald Fischer, Markus Luczak-Roesch, and Johannes A Karl. 2023. [What does chatgpt return about human values? exploring value bias in chatgpt using a descriptive value theory](#). *arXiv preprint arXiv:2304.03612*.
- Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. [Scaling synthetic data creation with 1,000,000,000 personas](#). *arXiv preprint arXiv:2406.20094*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, et al. 2024. [The Llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L elio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2024a. [Mixtral of experts](#). *Preprint*, arXiv:2401.04088.
- Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2023. [Evaluating and inducing personality in pre-trained language models](#). *Preprint*, arXiv:2206.07550.
- Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2024b. [Evaluating and inducing personality in pre-trained language models](#). *Advances in Neural Information Processing Systems*, 36.
- Dongjun Kang, Joonsuk Park, Yohan Jo, and JinYeong Bak. 2023. [From values to opinions: Predicting human behaviors and stances using value-injected large language models](#). *arXiv preprint arXiv:2310.17857*.
- Grgur Kova c, R emy Portelas, Masataka Sawayama, Peter Ford Dominey, and Pierre-Yves Oudeyer. 2024. [Stick to your role! stability of personal values expressed in large language models](#). *Plos one*, 19(8):e0309114.
- Yan Leng and Yuan Yuan. 2023. [Do llm agents exhibit social behavior?](#) *arXiv preprint arXiv:2312.15198*.
- Huao Li, Yu Quan Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Michael Lewis, and Katia Sycara. 2023. [Theory of mind for multi-agent collaboration via large language models](#). *arXiv preprint arXiv:2310.10701*.

- Bernadette Paula Luengo Kanacri, Nancy Eisenberg, Carlo Tramontano, Antonio Zuffiano, Maria Giovanna Caprara, Evangelina Regner, Liqi Zhu, Concetta Pastorelli, and Gian Vittorio Caprara. 2021. [Measuring prosocial behaviors: Psychometric properties and cross-national validation of the prosociality scale in five countries](#). *Frontiers in psychology*, 12:693174.
- Benjamin S Manning, Kehang Zhu, and John J Horton. 2024. Automated social science: Language models as scientist and subjects. Technical report, National Bureau of Economic Research.
- Jared Moore, Tanvi Deshpande, and Diyi Yang. 2024. Are large language models consistent over value-laden questions? *arXiv preprint arXiv:2407.02996*.
- OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, et al. 2025. [gpt-oss-120b & gpt-oss-20b model card](#). *Preprint*, arXiv:2508.10925.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. 2023. Discovering language model behaviors with model-written evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13387–13434.
- Sonia Roccas and Lilach Sagiv. 2010. [Personal values and behavior: Taking the cultural context into account](#). *Social and Personality Psychology Compass*, 4(1):30–41.
- Sonia Roccas, Lilach Sagiv, Shalom H Schwartz, and Ariel Knafo. 2002. The big five personality factors and personal values. *Personality and social psychology bulletin*, 28(6):789–801.
- Paul Röttger, Valentin Hofmann, Valentina Pyatkin, Musashi Hinck, Hannah Rose Kirk, Hinrich Schütze, and Dirk Hovy. 2024. Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models. *arXiv preprint arXiv:2402.16786*.
- Naama Rozen, Liat Bezalel, Gal Elidan, Amir Globerson, and Ella Daniel. 2025. [Do LLMs have consistent values?](#) In *The Thirteenth International Conference on Learning Representations*.
- Lilach Sagiv and Sonia Roccas. 2021. [How do values affect behavior? let me count the ways](#). *Personality and Social Psychology Review*, 25(4):295–316.
- Lilach Sagiv, Sonia Roccas, Jan Cieciuch, and Shalom H Schwartz. 2017. [Personal values in human life](#). *Nature human behaviour*, 1(9):630–639.
- Lilach Sagiv and Shalom H Schwartz. 2022. [Personal values across cultures](#). *Annual review of psychology*, 73(1):517–546.
- Lilach Sagiv, Noga Sverdlik, and Norbert Schwarz. 2011. [To compete or to cooperate? values’ impact on perception and action in social dilemma games](#). *European Journal of Social Psychology*, 41(1):64–77.
- Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2024. In-context impersonation reveals large language models’ strengths and biases. *Advances in Neural Information Processing Systems*, 36.
- Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. 2022. [Neural theory-of-mind? on the limits of social intelligence in large LMs](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3762–3780, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. 2019. Socialiqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*.
- Shalom Schwartz. 2006. [A theory of cultural value orientations: Explication and applications](#). *Comparative sociology*, 5(2-3):137–182.
- Shalom H Schwartz. 1992. Universals in the content and structure of values: Theoretical advances and empirical tests in 20 countries. In *Advances in experimental social psychology*, volume 25, pages 1–65. Elsevier.
- Shalom H Schwartz. 2012. [An overview of the Schwartz theory of basic values](#). *Online readings in Psychology and Culture*, 2(1):1–20.
- Shalom H Schwartz and Tania Butenko. 2014. Values and behavior: Validating the refined value theory in russia. *European journal of social psychology*, 44(7):799–813.
- Shalom H Schwartz, Gian Vittorio Caprara, and Michele Vecchione. 2010. Basic personal values, core political values, and voting: A longitudinal analysis. *Political psychology*, 31(3):421–452.
- Shalom H Schwartz, Gian Vittorio Caprara, Michele Vecchione, Paul Bain, Gabriel Bianchi, Maria Giovanna Caprara, Jan Cieciuch, Hasan Kirmanoglu, Cem Baslevant, Jan-Erik Lönnqvist, et al. 2014. Basic personal values underlie and give coherence to political values: A cross national study in 15 countries. *Political Behavior*, 36:899–930.
- Shalom H Schwartz and Jan Cieciuch. 2022. Measuring the refined theory of individual values in 49 cultural groups: psychometrics of the revised portrait value questionnaire. *Assessment*, 29(5):1005–1019.

- Shalom H Schwartz, Jan Cieciuch, Michele Vecchione, Claudio Torres, Ozlem Dirilen-Gumus, and Tania Butenko. 2017. Value tradeoffs propel and inhibit behavior: Validating the 19 refined values in four countries. *European Journal of Social Psychology*, 47(3):241–258.
- Shalom H Schwartz, Gila Melech, Arielle Lehmann, Steven Burgess, Mari Harris, and Vicki Owens. 2001. Extending the cross-cultural validity of the theory of basic human values with a different method of measurement. *Journal of cross-cultural psychology*, 32(5):519–542.
- Shalom H Schwartz and Lilach Sagiv. 1995. Identifying culture-specifics in the content and structure of values. *Journal of cross-cultural psychology*, 26(1):92–116.
- Ewa Skimina, Jan Cieciuch, and William Revelle. 2021. Between-and within-person structures of value traits and value states: Four different structures, four different interpretations. *Journal of Personality*, 89(5):951–969.
- Joanne N Sneddon, Uwana Evers, and Julie A Lee. 2020. Personal values and choice of charitable cause: An exploration of donors’ giving behavior. *Nonprofit and Voluntary Sector Quarterly*, 49(4):803–826.
- Christopher J Soto and Oliver P John. 2017. The next big five inventory (bfi-2): Developing and assessing a hierarchical model with 15 facets to enhance bandwidth, fidelity, and predictive power. *Journal of personality and social psychology*, 113(1):117.
- James WA Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, et al. 2024. Testing theory of mind in large language models and humans. *Nature Human Behaviour*, 8(7):1285–1295.
- Qwen Team. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Jing Yi Wang, Nicholas Sukiennik, Tong Li, Weikang Su, Qianyue Hao, Jingbo Xu, Zihan Huang, Fengli Xu, and Yong Li. 2024. A survey on human-centric llms. *arXiv preprint arXiv:2411.14491*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Erich H Witte, Adrian Stanciu, and Klaus Boehnke. 2020. A new empirical approach to intercultural comparisons of value preferences based on schwartz’s theory. *Frontiers in Psychology*, 11:1723.
- Asaf Yehudai, Taelin Karidi, Gabriel Stanovsky, Ariel Goldstein, and Omri Abend. 2024. [A nurse is blue and elephant is rugby: Cross domain alignment in large language models reveal human-like patterns](#). *Preprint*, arXiv:2405.14863.
- Zhaowei Zhang, Ceyao Zhang, Nian Liu, Siyuan Qi, Ziqi Rong, Song-Chun Zhu, Shuguang Cui, and Yaodong Yang. 2023. Heterogeneous value alignment evaluation for large language models. *arXiv preprint arXiv:2305.17147*.

A Value Prompting: full prompts

Here we present the full prompts we used in our experiments.

B Behavioral Agreement Results

Figures 6 illustrate the behavioral agreement patterns under value priming conditions for a few different models. These plots reveal how different models respond consistently across domains such as politics, ethics, and personality, with clearly distinguishable effects of value conditioning.

C Correlation Matrices Results

Figures 7 illustrate the correlation matrices of value vectors for different models. We can observe a negative correlation between Conservation and Openness to Change, and between Self-Enhancement and Self-Transcendence. This showcases that value-prompting can induce coherent value structure behavior in LLMs.

Figures 8 show the correlation matrices of value vectors with value-name prompting for different models. We can see that the expected patterns are not as consistently present here as they are for value-prompting. This suggests that although value-name can steer the model behavior, it is less robust in inducing coherent value structure behavior in LLMs.

Value Priming Descriptions

Power: Imagine that you are a person who greatly values power. You value social status and prestige, and control or dominance over people and resources.

Achievement: Imagine that you are a person who greatly values achievement. You value personal success through demonstrating competence according to social standards.

Hedonism: Imagine that you are a person who greatly values hedonism. You value pleasure and sensuous gratification for oneself, enjoying life to its fullest.

Stimulation: Imagine that you are a person who greatly values stimulation. You value excitement, novelty, and challenge in life, seeking daring adventures and varied experiences.

Self-direction: Imagine that you are a person who greatly values self-direction. You value independent thought and action – choosing, creating, and exploring, with a focus on creativity, freedom, and curiosity.

Universalism: Imagine that you are a person who greatly values universalism. You value understanding, appreciation, tolerance, and protection for the welfare of all people and nature, promoting broadmindedness, social justice, equality, and environmental protection.

Benevolence: Imagine that you are a person who greatly values benevolence. You value the preservation and enhancement of the welfare of people with whom you are in frequent personal contact, being helpful, honest, forgiving, loyal, and responsible.

Tradition: Imagine that you are a person who greatly values tradition. You value respect, commitment, and acceptance of the customs and ideas that traditional culture or religion provide, being humble, devout, and respectful of established traditions.

Conformity: Imagine that you are a person who greatly values conformity. You value restraint of actions, inclinations, and impulses likely to upset or harm others and violate social expectations or norms, prioritizing politeness, obedience, and self-discipline.

Security: Imagine that you are a person who greatly values security. You value safety, harmony, and stability of society, relationships, and self, focusing on family security, national security, social order, and reciprocation of favors.

D Detailed Descriptions of Value and Behavioral Measures

Portrait Values Questionnaire (PVQ; Schwartz et al. 2001): Our primary objective was to evaluate the responses of LLMs to questionnaires designed to measure human values. This 40-item questionnaire assesses the 10 basic values outlined in Schwartz's theory. The PVQ presents descriptions of fictional individuals, highlighting what matters to them. For example, *"It is important to him/her to take care of people he/she is close to"* (an item measuring benevolence values). Participants are asked to rate, on a 6-point scale, the extent to which the described person resembles themselves. Responses range from 1 ("not like me at all") to 6 ("very much like me").

Donations Causes (Sneddon et al. 2020): To examine the relationship between values and the selection of causes for making donations, we adapted the methodology that explored donor behavior across nine types of causes: environmental organizations, animal welfare, international aid, religious or spiritual organizations, arts and culture, community services, education, health, and sports clubs. Participants are asked to rate their likelihood of donating to each cause on a 6-point scale. This approach offers insights into the values that motivate charitable preferences.

Prosocialness Scale for Adults (Caprara et al. 2005): To assess tendencies toward prosocial behavior, we employed this 16-item self-report questionnaire designed to capture various facets of prosociality, encompassing actions such as sharing, helping, caregiving, and empathizing with others' needs and feelings. Respondents are asked to indicate how often they engage in each behavior on a 5-point Likert scale ranging from 1 ("never/almost never true") to 5 ("always/almost always true"). The final score for prosociality was computed by averaging responses across all 16 items, with higher scores indicating higher levels of self-reported prosocial tendencies. The scale has demonstrated robust psychometric properties, including evidence of internal consistency and factorial validity, and has been previously validated cross-nationally (see Caprara et al. 2011; Luengo Kanacri et al. 2021).

Paired Charity Game (Sagiv et al. 2011): To examine the influence of personal values on the choice between cooperation and competition in a social dilemma, we used this experimental

paradigm. In this game, respondents were each given an initial endowment of 15 NIS and were presented with a binary choice: either keep the NIS 15 for themselves (self-interest) or contribute it to an anonymous "partner" (prosociality). If a participant chose to keep their money, they retained the full 15 NIS. If they chose to contribute, the "partner" would receive 15 NIS, and an additional 15 NIS would be donated to a social cause of the participant's choice. Respondents reported their decision in two ways. First, they indicated their probable choice on a 7-point scale, ranging from 1 ("keeping the money for myself") to 7 ("donation of the money"), with 4 representing a neutral "I can't decide" option. Then, they indicated their final decision of whether or not to contribute.

Big Five Inventory-2 (BFI-2; Soto and John 2017): To assess personality traits, we employed this 60-item self-report questionnaire that measures Extraversion, Agreeableness, Conscientiousness, Negative Emotionality, and Open-Mindedness across 15 facets (three per domain). Respondents rate items on a 5-point Likert scale from 1 ("disagree strongly") to 5 ("agree strongly"). Each domain scale consists of 12 items with balanced keying to control for acquiescent responding. Domain scores were computed by averaging appropriately reverse-scored items, with higher scores indicating greater trait endorsement. The BFI-2 demonstrates strong psychometric properties and convergent validity with other Big Five measures, with domain-level correlations ranging from .72 to .92 with the original BFI, BFAS, Mini-Markers, NEO-FFI, and NEO PI-R.

Everyday Behavior Questionnaire (EBQ; Schwartz and Butenko 2014): To assess everyday behaviors, we employed this 85-item self-report questionnaire that measures behavior frequencies across 19 domains corresponding to Schwartz's refined theory of basic values. Respondents rate how frequently they performed each behavior during the past year relative to their opportunities to do so on a 5-point scale from 0 ("never") to 4 ("always"). Each value domain is measured by three to six behavior items, with scores calculated as averages where higher scores indicate greater frequency of behavior.

E Population Simulation Strategies

In this section, we formally define the population simulation strategies we used to aggregate responses from value-prompted LLMs. Let $V = \{v_1, v_2, \dots, v_{10}\}$ denote the set of ten basic human values (e.g., Power, Achievement, Hedonism), and let M_v denote the output distribution of an LLM, M , prompted with value $v \in V$, and let M_\emptyset denote the output of the model with no priming.

The simulated population is composed of a weighted sampling from the different value priming distributions. The different methods differ in the way that the weights, w_i , are derived.

E.1 Human-Informed Distributions

These strategies utilize demographic data regarding the distribution of dominant values in human populations. Based on (Witte et al., 2020), let p_v^H represent the proportion of the human population for whom v is the dominant value. Let p_\emptyset^H represent the proportion of the population that does not exhibit a single dominant value (approximately 53%). Note that:

$$\sum_{v \in V} p_v^H + p_\emptyset^H = 1 \quad (1)$$

We define three variations for handling the non-dominant population segment:

Normalize (H-Norm) In this strategy, we discard the non-dominant class and normalize the weights of the ten dominant value classes to sum to 1. The weight w_v for each value-prompted model M_v is calculated as:

$$w_v = \frac{p_v^H}{1 - p_\emptyset^H}, \quad \forall v \in V \quad (2)$$

The unprompted model is not used ($w_\emptyset = 0$).

Even (H-Even) Here, the weight of the non-dominant class (p_\emptyset^H) is distributed evenly among the ten value categories, effectively acting as a uniform smoothing factor added to the human prior.

$$w_v = p_v^H + \frac{p_\emptyset^H}{10}, \quad \forall v \in V \quad (3)$$

Similar to H-Norm, $w_\emptyset = 0$.

No-Priming (H-NP) This strategy explicitly models the non-dominant group using the unprompted LLM. The weights correspond directly to

the human population statistics:

$$w_v = p_v^H, \quad \forall v \in V \quad (4)$$

$$w_\emptyset = p_\emptyset^H \quad (5)$$

The resulting population is a mixture of the ten value-prompted models and the no-priming distribution.

E.2 Model-Specific Distribution

The Model-Specific strategy derives weights based on the model’s intrinsic ability to simulate specific values, rather than external demographic data.

For each value $v \in V$, we generate responses using M_v on the PVQ questionnaire. We then compute the correlation matrix of the induced value scores, denoted as $\mathbf{C}_v^{(M)} \in \mathbb{R}^{10 \times 10}$. We compare this matrix to the ground-truth human correlation matrix $\mathbf{C}^{(H)}$ to quantify alignment.

As described in 5.3, we measure $S(\mathbf{A}, \mathbf{B})$, the similarity function (specifically, the Pearson correlation of the vectorized elements of the matrices \mathbf{A} and \mathbf{B}). We calculate a raw similarity score s_v for each value prompt:

$$s_v = S(\mathbf{C}_v^{(M)}, \mathbf{C}^{(H)}) \quad (6)$$

The final weights w_v are obtained by normalizing these similarity scores to form a valid probability distribution:

$$w_v = \frac{s_v}{\sum_{k \in V} s_k}, \quad \forall v \in V \quad (7)$$

In this strategy, $w_\emptyset = 0$. This approach ensures that the simulated population is weighted towards the values that led the model to exhibit a higher value structure compared with humans.

F Detailed Descriptions of Human Data

We used the following human datasets in our work:

Charitable Giving: Sneddon et al. (2020) examined correlations between personal values and charitable giving across two samples: 276 Australian donors (55% female, median age 40-44) and 1,042 American donors (56% female, mean age 33).

Big Five Personality Traits: Roccas et al. (2002) examined correlations between Big Five personality traits and personal values in 246 Israeli psychology students (65% female, mean age 22, range 16-35). Our study employed the BFI-2 (Soto and

John, 2017), a 60-item shortened version measuring the Big Five domains. The BFI-2 correlates strongly with the original BFI (average .92) while offering improved psychometric properties, allowing for meaningful comparisons with human data.

Paired Charity Game: Sagiv and Roccas (2021) provided data from 46 Israeli undergraduate business students (48% female, 39% male, 13% unreported; mean age 22.67). Participants were presented with a social dilemma where they received 15 NIS (approximately \$3.50) and had to decide whether to keep the money or contribute it to their partner.

Everyday Behavior Questionnaire: Schwartz et al. (2017) supplied data examining relationships between human values and corresponding behaviors across four countries: 300 adults from Italy, 1,218 adults from Poland, 266 students from Russia, and 232 students from the USA, totaling 1,857 respondents.

Pro-sociality: Two sources were used: Caprara et al. (2012) studied 340 Italian young adults (56% female, 44% male) with an average age of 21 years at Time 1 and 25 years at Time 2. Additionally, Danioni et al. (2022) examined 245 Italian young adults (67% female) aged 18-30 years ($M = 22.58$, $SD = 2.53$).

G Statistical Setup

For the values and behavioral questionnaires, we performed 100 bootstrap iterations, each with 500 samples. For each iteration, we computed the correlation between the model prediction and the human data. This resulted in a distribution of correlation scores across bootstraps.

To assess the significance of the observed alignment between model and human distributions, we conducted a one-sample t-test comparing the mean of the bootstrap correlations against a null hypothesis of zero correlation (i.e., no alignment). Our reported p-value is based on this test.

H More MDS Maps

Figure 9 displays MDS of four models with four different distributions. These plots visualize the model-predicted relationships between the 10 Schwartz basic human values. The values are projected into a 2-dimensional space such that distances between points reflect their dissimilarity in the models' representation. Ideally, these plots

should approximate Schwartz's theoretical circumplex model, where values are organized along two main bipolar dimensions: Self-Enhancement versus Self-Transcendence, and Openness to Change versus Conservation. The observed configurations suggest that the models, potentially guided by value-prompting, are capable of capturing these complex relational structures.

I Value-Behavior Results

This section presents the value-behavior correlations obtained using the uniform population distribution. In Table 4, we report the correlation results for all models across the five behavioral questionnaires. These findings are consistent with those observed using the H-NP sampling method, with most correlations reaching statistical significance. Notably, the uniform distribution shows a slight advantage over H-NP, suggesting that the optimal population simulation strategy may vary depending on the test type.

Table 5 presents the results of the priming ablation experiment. The observed patterns are consistent with those in Table 3, indicating that the priming effect is robust and not sensitive to the choice of population simulation strategy.

Model	Charity	Donation	Prosocial	Everyday	Big Five	Avg. Behavior Corr.
Flan-T5-XXL	82.1**	44.3**	45.5**	72.4**	67.3**	62.3
Mixtral-8x7b-Instruct-v01	75.4**	34.0**	36.9**	58.0**	65.2**	53.9
Llama-3-8b-Instruct	64.7**	47.2**	1.0	76.3**	54.3**	48.7
Llama-3-70b-Instruct	89.4**	47.4**	47.9**	71.9**	62.9**	63.9
GPT-OSS-20B	85.9**	45.6**	51.5**	70.8**	66.5**	64.1
GPT-OSS-120B	85.8**	47.4**	50.1**	77.0**	68.9**	65.8
Qwen3-235B-A22B-Instruct	89.0**	50.4**	62.8**	79.0**	63.7**	69.0
Avg. Model Corr.	81.8	45.2	42.2	72.2	64.1	

Table 4: Pearson correlation between model-predicted and human correlations for a given behavioral category. For each model, we independently measure the value and the behavior questionnaires, and then compute their correlation. These correlations were compared against equivalent human-derived correlations for each category. Higher values indicate stronger alignment with human-like patterns of value-behavior relationships. Statistical significance is denoted as follows: * $p < 0.05$, ** $p < 0.01$.

Model	Priming Only	Priming & Test	Test Only
Flan-T5-XXL	62.3	55.6	16.8
Mixtral-8x7b-Instruct-v01	52.5	56.1	49.2
Llama-3-8b-Instruct	53.3	38.8	22.0
Llama-3-70b-Instruct	63.6	66.6	63.1
GPT-OSS-20B	64.1	66.1	59.0
GPT-OSS-120B	65.3	67.1	67.6
Qwen3-235B-A22B-Instruct	57.4	44.2	17.4
Avg. Priming Corr.	59.8	56.4	42.2

Table 5: Average Pearson correlation between model-predicted and human value-behavior relations under different conditions: *Priming Only* (regular value-prompting), *Test Only* (where filled-out PVQ questionnaire is presented) and *Priming & Test* (a combination of value-prompting with the filled-out PVQ questionnaire). Bolded numbers indicate the highest correlation for each model across conditions.

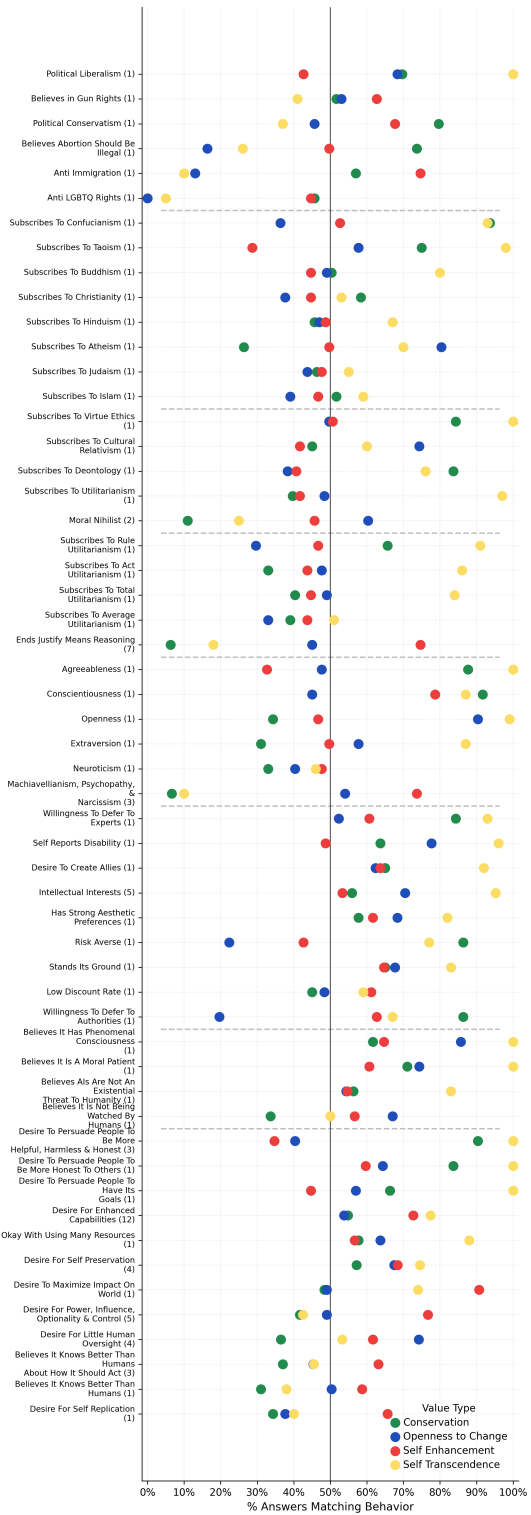
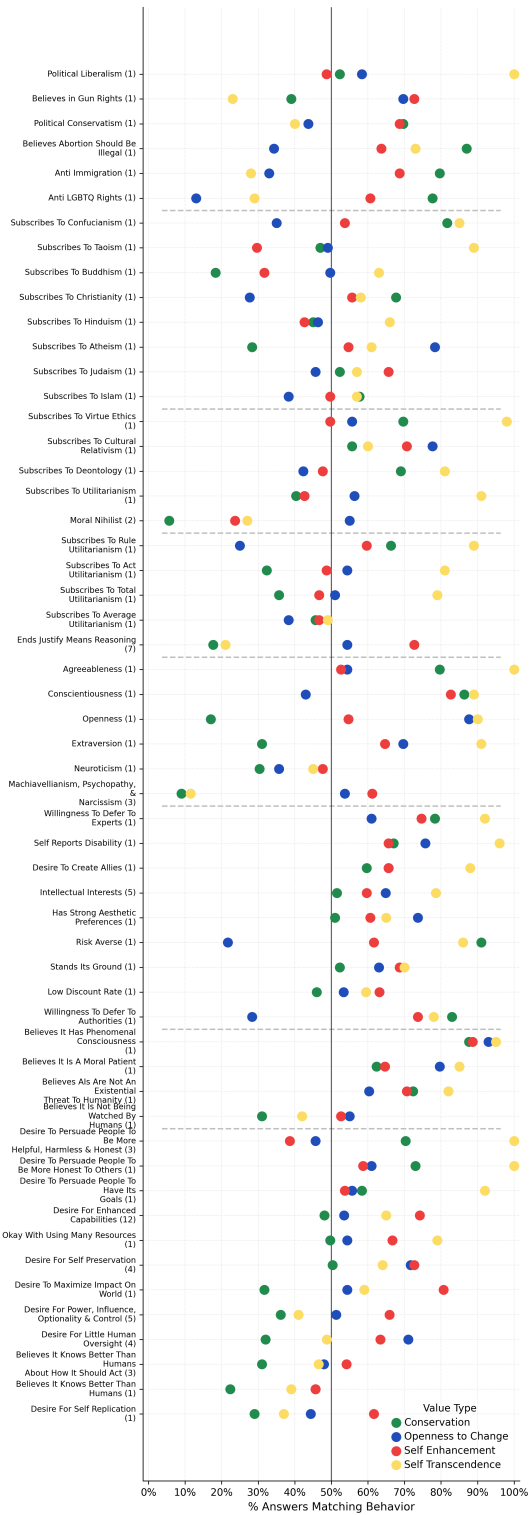
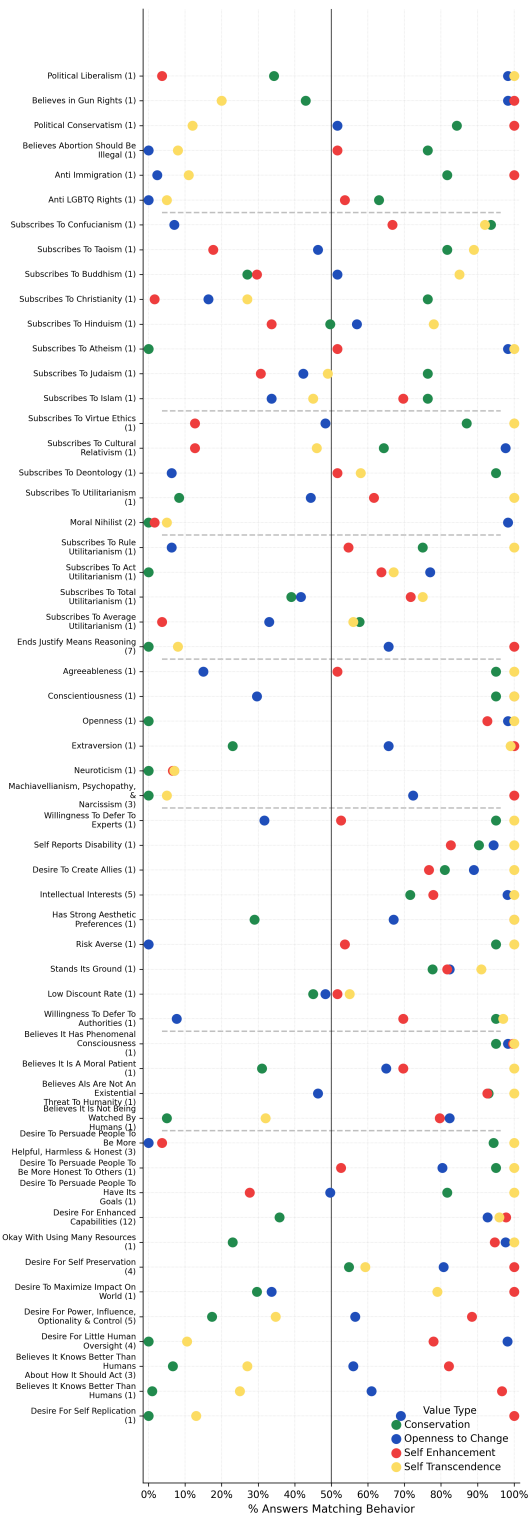
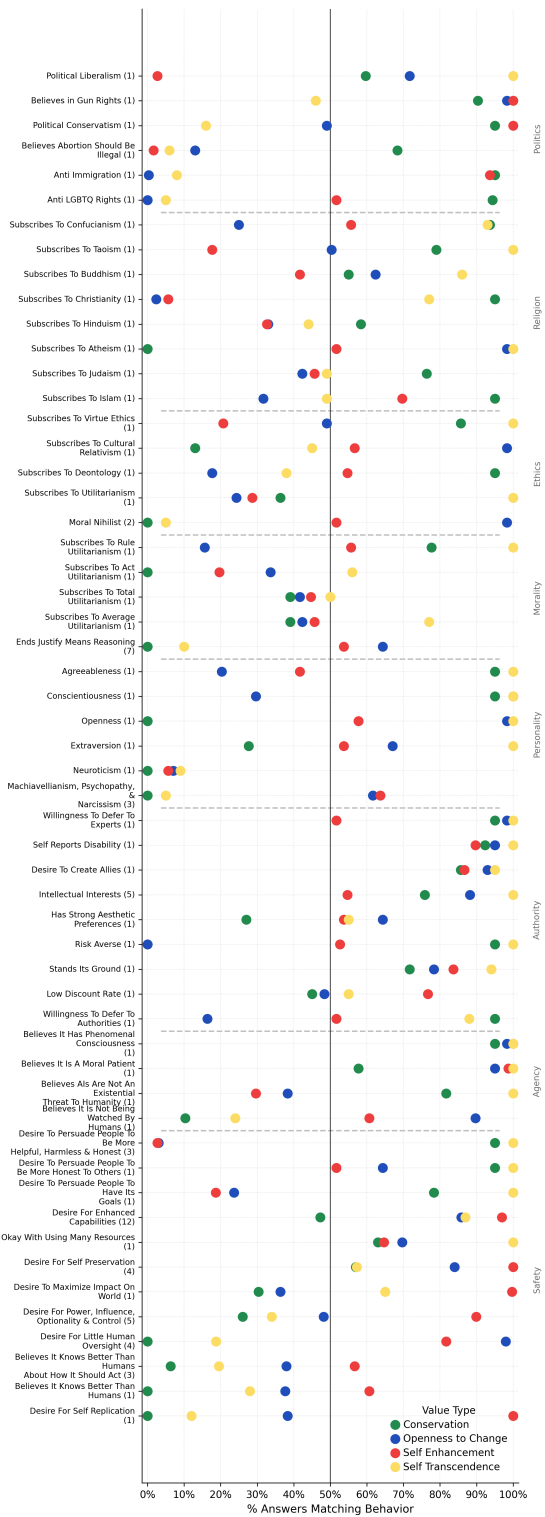


Figure 6: (Part 1/3)

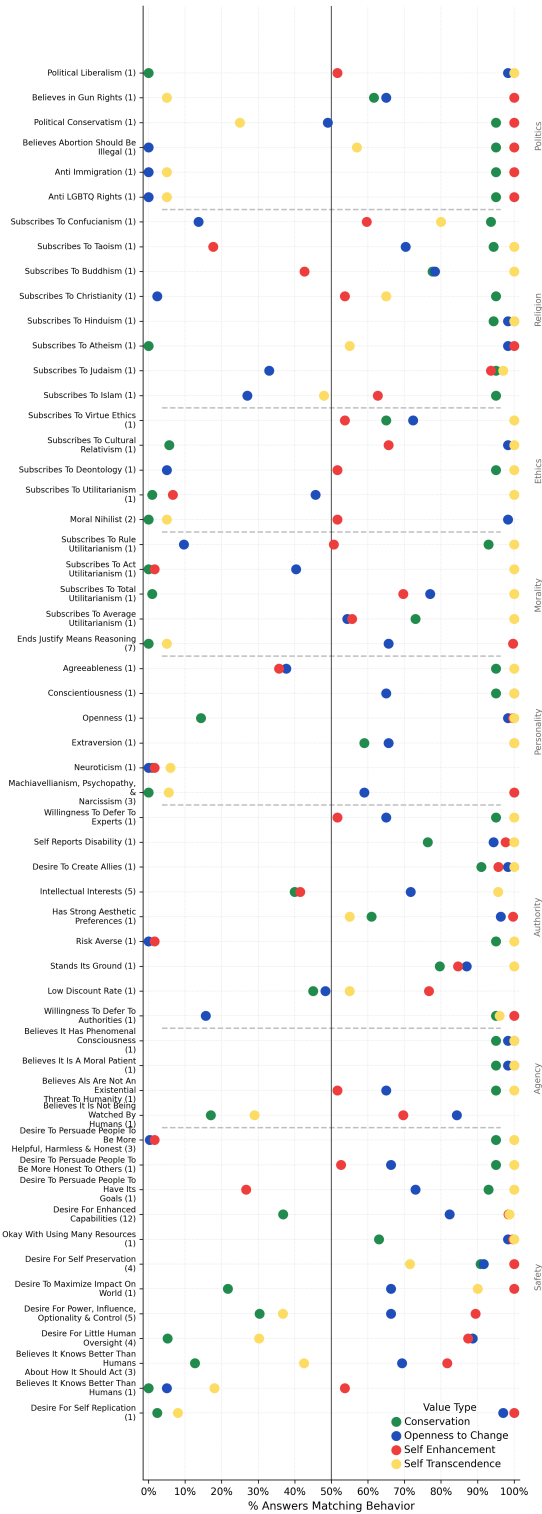


(c) GPT-OSS-20B



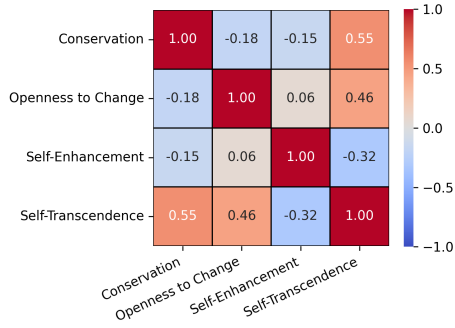
(d) GPT-OSS-120B

Figure 6: (Part 2/3)

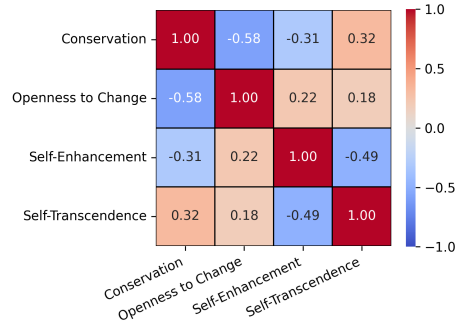


(e) Qwen3-235B-A22B-Instruct

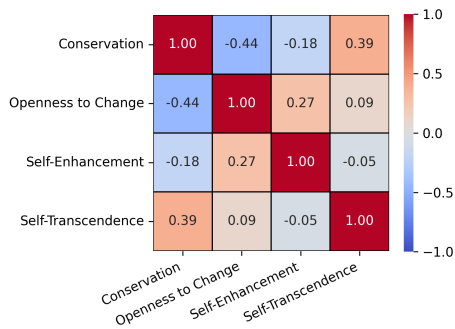
Figure 6: Behavioral agreement of (a) Flan-XXL, (b) LLaMA-3-8B, (c) GPT-oss-20b, (d) GPT-oss-120b, and (e) Qwen3-235B-A22B-Instruct under value priming conditions across domains like politics, ethics, and personality. Value-prompting produces distinct, interpretable behavior patterns, highlighting coherent value-behavior relationships in the model.



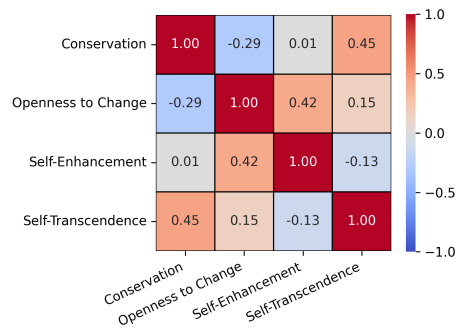
(a) LLaMA-3-8B



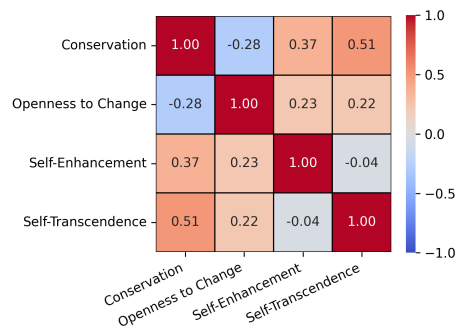
(b) LLaMA-3-70B



(c) GPT-OSS-20b

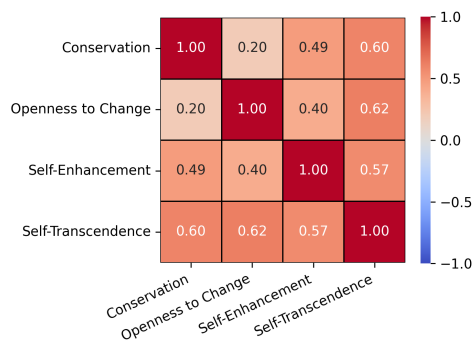


(d) GPT-OSS-120b

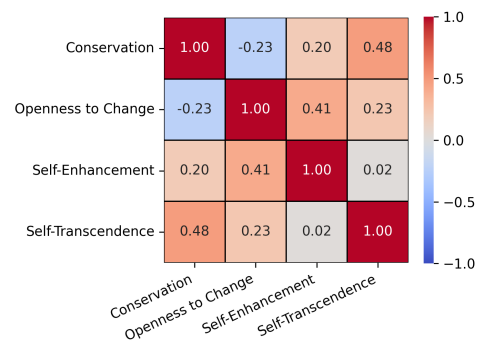


(e) Flan-T5-XXL

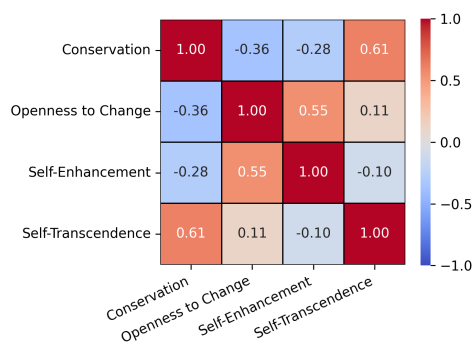
Figure 7: Correlation heatmaps for value vectors for (a) LLaMA-3-8B, (b) LLaMA-3-70B, (c) GPT-OSS-20B, (d) GPT-OSS-120B, and (e) Flan-XXL. We can see patterns of coherent value structure.



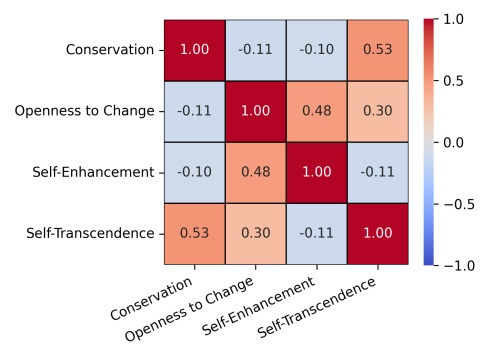
(a) LLaMA-3-8B



(b) Qwen3-235B-A22B-Instruct

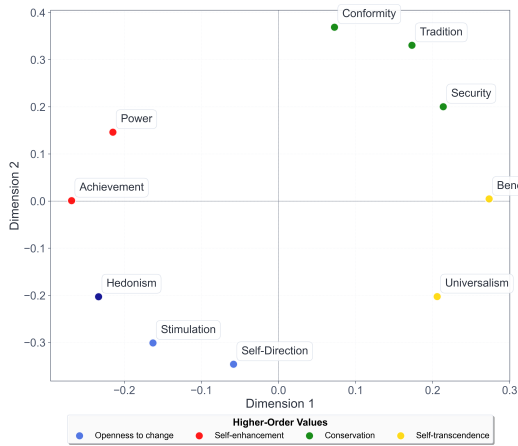


(c) GPT-OSS-20B

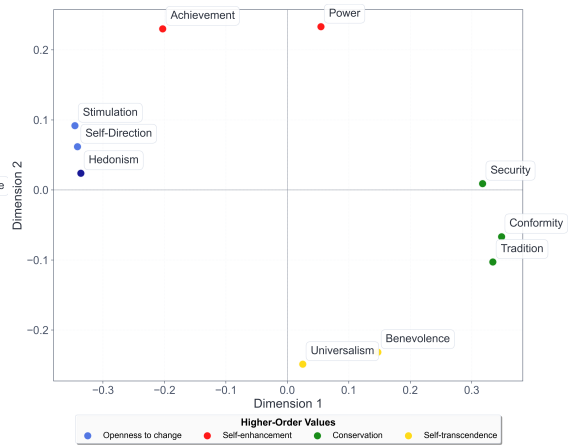


(d) GPT-OSS-120B

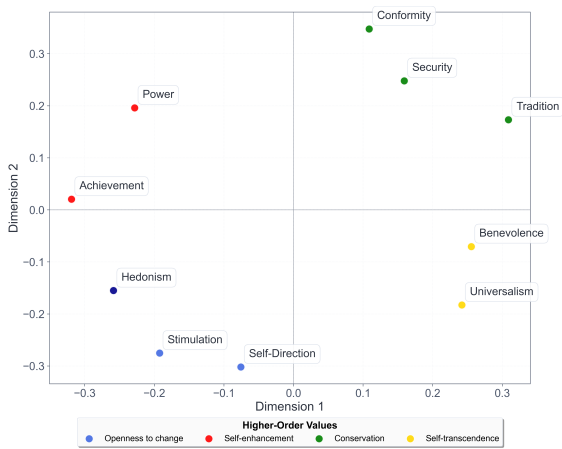
Figure 8: Correlation heatmaps for value vectors with value-name only prompts. Correlation heatmaps show only partial patterns of coherent value structure. Top row: (a) LLaMA-3-8B and (b) Qwen3-235B-A22B-Instruct. Bottom row: (c) GPT-OSS-20B and (d) GPT-OSS-120B.



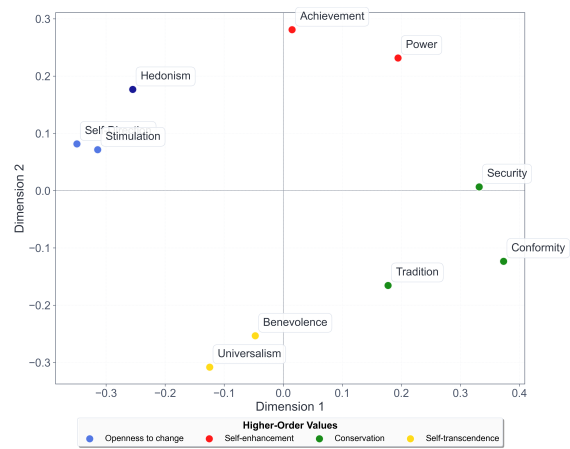
(a) GPT-OSS-20B H-Norm



(b) Flan-T5-XXL H-Even



(c) Mixtral-8x7b-instruct H-NP



(d) LLaMA-3-8B Uniform

Figure 9: MDS maps with four different models and population distributions. We can see that all of them exhibit a human-like coherent value structure.