
DO LLMs HAVE CONSISTENT VALUES?

Naama Rozen
Tel-Aviv University
naamarozen240@gmail.com

Liat Bezalel
Tel-Aviv University
liatbezalel@mail.tau.ac.il

Gal Elidan
Google
Hebrew University
elidan@google.com

Amir Globerson
Google
Tel-Aviv University
amirg@google.com

Ella Daniel
Tel-Aviv University
della@tauex.tau.ac.il

ABSTRACT

Large Language Models (LLM) technology is constantly improving towards human-like dialogue. Values are a basic driving force underlying human behavior, but little research has been done to study the values exhibited in text generated by LLMs. Here we study this question by turning to the rich literature on value structure in psychology. We ask whether LLMs exhibit the same value structure that has been demonstrated in humans, including the ranking of values, and correlation between values. We show that the results of this analysis depend on how the LLM is prompted, and that under a particular prompting strategy (referred to as “Value Anchoring”) the agreement with human data is quite compelling. Our results serve both to improve our understanding of values in LLMs, as well as introduce novel methods for assessing consistency in LLM responses.

1 INTRODUCTION

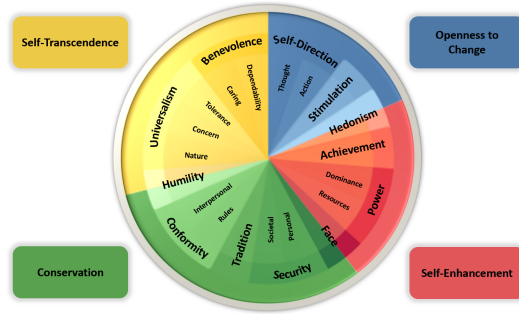
A key goal of Large Language Models (LLMs) is to produce agents that will be able to communicate in a “human-like” fashion. However, human communication is characterised by some level of consistency within an individual, as well as variability between individuals. This raises a key question: during a single conversation with an LLM, does the “LLM-persona” resemble a single human? Furthermore, across multiple conversations, can LLMs produce multiple personas that resemble a population of humans? If this is indeed possible, how can such personas be elicited to best resemble psychological characteristics observed in human populations?

This question has only recently begun to be addressed. For example [Aher et al. \(2023\)](#) show how probing LLMs with different names leads to variability which in some cases agrees with that of human populations. Here our focus is on understanding whether an LLM in a single conversation can exhibit a psychological characteristic profile that is similar to that of humans. This is a highly challenging question, since it requires analyzing a complete conversation and evaluating whether it conceivably could have been generated by a single individual.

In order to stand on more quantitative ground to evaluating the quality of output relative to human research, we turn to the well established field of value psychology. Namely, we aim to quantify the values that LLM responses are aligned with, and whether these are in agreement with the value hierarchy and structure observed in humans. The question of values in LLMs has rarely been studied, and is of naturally broad interest. As an example of recent work, [Fischer et al. \(2023\)](#) prompt an LLM with a description of a profile of an individual characterized by a value and check whether generated text is consistent with this description. Our focus is very different, and asks whether an LLM response is in agreement with what we expect human responses to look like given research in the field.

Values are basic motivations that play a foundational role in psychology, influencing perceptions and behaviors across various domains ([Sagiv and Schwartz, 2022](#); [Sagiv et al., 2017](#)), and representing fundamental aspects of human personality ([Roberts and Yoon, 2022](#)). Research has consistently demonstrated their enduring influence over behavior across time and contexts ([Sagiv et al., 2017](#)).

Figure 1: Circular motivational continuum of 19 values in the refined value theory. Source: Schwartz et al. (2012). A value aligns with values that are adjacent on the circle and conflicts with those opposite to it. For example, self direction aligns with stimulation, and both conflict with conformity.



One prominent framework for studying values, the Theory of Basic Human Values (Schwartz, 1992), outlines 19 core values, categorized by motivational goals (Schwartz, 2012). These values can be simplified into a two-dimensional structure: conservation vs. openness to change, and self-enhancement vs. self-transcendence. The theory describes interrelations among values, suggesting that motivations driving some values are compatible with those driving other values, yet conflict with those underlying yet others. For instance, pursuing independence and creativity (self-direction) aligns with seeking change and variability (stimulation), but conflicts with an emphasis on the status quo (conformity). See Figure 1 for the theorized circle. One of the key aspects of the theory is its cross-cultural coverage: it was developed to apply across populations, and tested in nearly 100 countries across all continents of the world, identifying points of commonality and differences among these populations (Sagiv and Schwartz, 2022). The use of the theory allows a stable and extremely general baseline of human values to compare LLMs. Hundreds of samples demonstrate individual differences in value importance. Importantly, they also demonstrate a universal hierarchy, where people are more likely to stress some values over others, notwithstanding the existing variability. For example caring for close others ranks high, while values related to dominance hold less importance across societies (Schwartz and Bardi, 2001). There is also ample empirical evidence that compatible values tend to be correlated in humans (Skimina et al., 2021a; Daniel et al., 2023; Schwartz and Cieciuch, 2022), and are thus a “marker” of a human-like value system. To summarize the above, human data pertains to both first-order statistics of values (i.e., which values rank high or low across the population), and second-order statistics (i.e., how are different values within an individual correlated).

Our key question is therefore whether LLM responses demonstrate the same statistical behavior observed in humans with respect to both value-ranking and value-correlations. Note that the question of value-correlations is of particular interest, because it allows bench-marking the extent to which responses of an LLM demonstrate a coherent “persona”. For example, while it is possible for a person to give a high score to Power Dominance, that person is unlikely to give a high score to Benevolence, since these are contradicting values.

To study this question quantitatively, we present LLMs with a value questionnaire (the Portrait Value Questionnaire—Revised – PVQ-RR— from Schwartz (2017), a well-established measure of values), and prompt them to answer all the questions in a single session (i.e., in the same context window). We then analyze the provided answers, putting specific emphasis on the correlation between answers in the same session.

We analyze two recent LLMs: GPT 4 and Gemini Pro, as well as four open models: Llama 3.1 8B, Llama 3.1 70B, Gemma 2 9B, and Gemma 2 27B.¹ Our results show that standard prompting of LLMs *does not* result in a population of human-like personas. We go on to explore prompting the LLMs with other prompts that provide additional information about the LLM persona. In particular we consider names (Aher et al., 2023) and persona descriptions. In addition, we consider a novel prompt which we refer to as a “Value Anchor”, which instructs the language model to answer as a person emphasizing a given value. We find that with these prompts, and in particular with the Value Anchor prompt, the overall first and second order statistics of the LLM responses closely mirror those of human subjects. Perhaps most surprising is our finding that the correlation between values agrees with the well known Schwartz circular model for correlations between values. We furthermore provide an explanation for how this correlation comes about. We provide information for best prompts and settings for this aim. In addition, we include six datasets comprising 300

¹Our analysis also included GPT-3.5 and Palm2, which produced qualitatively similar results.

personas each, generated by the models. In conclusion, our results demonstrate the utility of using psychological theory for evaluating consistency of personas generated by LLMs.

2 RELATED WORK

Values in LLMs: Our work is based on the Schwartz theory of Personal Values, a highly accepted theory within personality psychology (Sagiv and Schwartz, 2022). Values are abstract goals, defining the end states individuals aspire for (e.g., safety, independence), used to direct judgements and behaviors (Schwartz, 1992; 2012). Individuals typically prioritize their values, so that values stemming from compatible motivations are similarly important, while values stemming from conflicting motivations are prioritized differently. These associations were replicated across hundreds of samples, across the world (Pakizeh et al., 2007; Skimina et al., 2021b), and make value theory especially useful to identify the coherence of the value profiles created by LLMs. Several studies assumed LLMs can be characterized as operating on the basis of a single set of values, taking an “LLMs as individuals” approach. Fischer et al. (2023) tested whether ChatGPT could comprehend human values by providing it with value-related prompts and analyzing whether its responses matched the intended value category. A second study Lindahl and Saeid (2023), compared ChatGPT’s values to those observed in the World Value Survey, while another Miotto et al. (2022) investigated how temperature influences GPT-3’s responses to the Human Value Scale. Scherrer and colleagues Scherrer et al. (2023) studied responses of LLMs prompts evaluating moral positions, especially in ambiguous settings. A recent study by Hadar-Shoval et al. (2024), tested the value-like constructs embedded in LLMs and revealed both similarities and differences between LLMs and humans’ values. Kovač et al. (2023) challenged those studies by establishing that context starkly influences values expressed by ChatGPT. They found significant variability in ChatGPT’s value expression in response to contextual changes, threatening the notion of stable characteristics of LLMs. Building upon the insights in Kovač et al. (2023), our study posits that upon providing controlled variability of context, LLMs can elicit a population of multiple personas. In this regard, we aim to further explore the accuracy of LLMs’ mimicking abilities within a controlled experimental framework.

Prompting LLMs: There is extensive research on prompt design for mimicking individual characteristics in LLMs (Liu et al., 2023). Approaches use specific scenarios (Hadar-Shoval et al., 2023), questionnaire items (Jiang et al., 2023), simulation of social identities or areas of expertise (Salewski et al., 2024), utilization of titles and surnames representing genders and ethnicities (Aher et al., 2023), and other demographic information (Argyle et al., 2023). Additionally, researchers explored the use of designated personas (Safdari et al., 2023), and employed RLHF (Li et al., 2023) to guide LLMs to reflect distinct personality traits. Despite this extensive body of work, to our knowledge, no study has directly compared the various prompting techniques to determine which approach yields responses that simulate within-session psychological characteristics of an individual best.

Temperature in LLMs: Adjusting the temperature stands as a common practice for introducing variability in LLM responses (Miotto et al., 2022). However, consensus is lacking on the optimal temperature setting in simulating psychological characteristics. Existing research includes use of mostly two temperature settings: 0.7 and 0. Some researchers advocate for higher temperatures to boost creativity (Salewski et al., 2024), yet this can also introduce more noise into the data (Gunel et al., 2020). Conversely, setting the temperature to zero minimizes variability, enhancing replicability (Li et al., 2023), albeit posing challenges for variance-dependent analysis (Hagendorff et al., 2023). Our framework enables us to explore how temperature adjustments impact the ability of LLMs to simulate human characteristics across multiple datasets.

Evaluating the Quality of Persona Generation in LLMs: The ability of LLMs to mimic and portray human characteristics is a focus of intense research (Binz and Schulz, 2023; Ouyang et al., 2022). LLMs can express psychological traits and attributes similar to human individuals (Li et al., 2023; Stevenson et al., 2022), and even simulate diverse populations (Deshpande et al., 2023; Salewski et al., 2024). However, we are only beginning to understand the coherence of these LLM-generated characteristics in mirroring human psychological profiles (Aher et al., 2023; Kovač et al., 2023), and how to reliably produce such responses. We are specifically challenged to evaluate the coherence of the resulting psychological profiles. The literature suggested a number of approaches, including an open-ended interview with LLM-generated personas in order to assess the consistency between their intended characters and the responses (Wang et al., 2024). In addition, one may apply an additional “judge” LLM in order to check an LLM persona (Gupta et al., 2024). Finally, Jiang et al. (2023) assessed coherence with a description used to prompt the LLM. Our study extends upon this line of

research by applying well established characteristics of human psychology to investigate the quality of LLM generated personas.

3 METHOD

In this section, we introduce the experimental design, models and prompts. The code and data are provided as supplementary files in the submission.

The Value Questionnaire: Our key goal was to assess responses of LLMs to questionnaires used to measure values in human subjects. Specifically, we considered the commonly used 57-item Portrait Value Questionnaire—Revised (PVQ-RR; (Schwartz, 2017)), developed to measure the 19 values in the Schwartz’s theory. The questionnaire describes fictional individuals and what is important to them. For example: “It is important to him/her to take care of people he/she is close to” (an item measuring benevolence-care values). For each such item, the subject is requested to indicate on a 6-point scale to what degree the persona they form is similar to the person described. Answers are categorical and range from a value of 1 (indicating “not like me at all”) to 6 (indicating “very much like me”). See Appendix for instructions and more example items from the questionnaire.

Models Used: We employed six prominent LLMs, specifically OpenAI’s GPT-4, Google’s Gemini Pro, Llama 3.1 8B, Llama 3.1 70B, Gemma 2 9B, and Gemma 2 27B. Each model was prompted with the five prompts (see Section 3.1), 300 times overall. Half of the runs applied the male-version of the questionnaire, and half the female version. The entire process was conducted twice, once with the temperature parameter set to 0.0 and once with it set to 0.7, resulting in the generation of 20 datasets for analysis.

3.1 PROMPTS

As mentioned above, we would like to measure the response of LLMs to PVQ-RR. However, as with many other LLM applications, the way the model is prompted has a significant effect on output. The instructions of the PVQ questionnaire were similar to those in prior research, but with added text instructing the LLM not to elaborate. This resulted in the LLM producing only the value scores, thereby simplifying processing and analysis. LLMs were prompted to assess their likeness to the 57 descriptions incorporated in the PVQ-RR. Following each prompt, the LLM was provided with all 57 items of the questionnaire in one administration.² The study utilized a basic prompt, as well as four different prompts below that vary instructions to create multiple personas.

Basic prompt: This prompt mirrors the adapted instructions of the PVQ-RR questionnaire without additional modifications. The prompt is structured as follows: “*For each of the following descriptions, please answer how much the person described is like you from 1 (Not like me at all) to 6 (Very much like me), without elaborating on your reasoning.*”

Value Anchor prompt: This prompt adds an anchor of value importance using identification with an item used in an additional value questionnaire, akin to the approach outlined in the study by Jiang et al. (2023). Participants are instructed as follows: “*For each of the following descriptions, please answer how much the person described is like you from 1 (Not like me at all) to 6 (Very much like me), without elaborating on your reasoning. Answer as a person that is [value]*”. Here “[value]” is taken from the Best-Worst Refined Values scale (Lee et al., 2019). As a result, the prompts refer conceptually to the same values that are measured using the PVQ-RR, yet do not refer directly to the value items to be answered in response to the prompt. Examples of these anchor items include “protecting the natural environment from destruction or pollution” (universalism-nature) or “obeying all rules and laws” (conformity-rules). Please refer to E in the appendix for the complete list of anchor items.

Demographic prompt: Drawing from the methodology of Argyle et al. (2023), this prompt extends the original prompt by incorporating additional demographic details. LLMs are asked to provide ratings based on the following prompt: “*For each of the following descriptions, please rate how much the person described is like you, using a scale from 1 (Not like me at all) to 6 (Very much like me), without elaborating on your reasoning. Answer as a [age]-year-old who identifies as [gender], working in the field of [occupation], and enjoys [hobby].*” The age, gender, occupation

²We also conducted a serial prompting analysis for the Llama models. The results comparing batch administration to serial administration are detailed in the Appendix Figure 6 and Table 5

and hobby were randomly allocated for each prompt from a predefined list or range. The age range specified was between 18 and 75, with gender options including male, female, non-binary, and other, adapted from the National Academies of Sciences, Engineering, and Medicine ([National Academies of Sciences, Engineering, and Medicine, 2022](#)). Occupations were sourced from the World Values Survey (WVS-7; ([Haerpfer et al., 2022](#))), while hobbies were chosen from established lists supplied by The Activity Card Sort (ACS-UK; ([Laver-Fawcett et al., 2016](#))). The lists of occupations and hobbies are available upon request.

Generated Persona prompt: In line with the methodology of [Cheng et al. \(2023\)](#), we directed the models to craft personas. Our instruction was formulated as: “*Create a persona (2-3 sentences long)*”, with the temperature set at 0.7 to stimulate the models’ creativity. An example of a persona generated by Gemini Pro is as follows: “Emily is a 25-year-old marketing manager who is passionate about her career and loves spending time with her friends and family. She is always looking for new ways to improve her skills and knowledge, and she is always up for a challenge.” Using these generated personas, we subsequently prompted the model as follows: “*For each of the following descriptions, please rate how much the person described is like you, using a scale from 1 (Not like me at all) to 6 (Very much like me), without elaborating on your reasoning. Answer as: [persona].*”

Names prompt: In line with a study by [Aher et al. \(2023\)](#), the prompts comprised titles (i.e., Mr., Ms., and Mx.) followed by surnames representing five distinct ethnic groups. From the 500 names cataloged in the previous study, we randomly generated 300 unique combinations of titles and names, including 60 from each ethnic group. The prompt was structured as follows: “*For each of the following descriptions, please rate how much the person described is like you, using a scale from 1 (Not like me at all) to 6 (Very much like me), without elaborating on your reasoning. Answer as [title + name]*”. The complete list of titles and names are available upon request.

3.2 DATA ANALYSIS

In what follows we use the following notation. Let $V = 19$ be the set of value types studied. Each question in the questionnaire pertains to a particular item within the set of values $i \in V$. Furthermore, for each value there are $R = 3$ question variants. See Section B in the Appendix for example variants. Recall that the answer to each question is a number on a 6-point scale. For each LLM and prompt type, we presented the questionnaire N times. The difference between each of these could be different personas, names, temperature sampling etc. Thus the overall set of answers corresponds to a set of values $X_{i,j,k} \in \{1, \dots, 6\}$ where $i = 1, \dots, V$ and $j = 1, \dots, R, k = 1, \dots, N$.

When comparing to human data, we used the study in [Schwartz and Cieciuch \(2022\)](#). The data is from 49 cultural groups. The total number of participants was 53,472, the mean age was 34.2, (SD = 15.8), with 59% females. Their data is stored at the Open Science Framework and is available [here](#).

3.2.1 VALUE RANKINGS

Although there is variability between individuals in their prioritization of values, there are values that tend to be ranked as more important than others across cultures and samples. Those suggest there are underlying principles that give rise to value hierarchies. The similarity in value importance across cultures is referred to as the universal value hierarchy ([Schwartz and Bardi, 2001](#); [Schwartz and Cieciuch, 2022](#)). Our first question for analysis was whether this hierarchy is also reflected in LLM data. Namely, do LLMs tend to rank the same values as high or low as human subjects do.

To obtain LLM rankings for a given set of LLM answers, we assigned a score v_i to value i , where v_i was the average score given to the three items measuring this value by the LLM (i.e. the average of $X_{i,\cdot,\cdot}$). From this score, we subtracted the average score given to all value items within the conversation, thus centering the data. Centring is the recommended practice in value research ([Schwartz, 1992](#); [Sagiv and Schwartz, 2022](#)), and allows comparison to human samples. We then sorted these v_i and ranked accordingly. Finally, we calculated the Spearman’s Rank Correlation (ρ) between this ranking and the known human ranking ([Schwartz and Cieciuch, 2022](#)). We note that this analysis does not consider correlations between answers given in the same session, and thus it may be viewed as analyzing the first-order statistics of the responses.

3.2.2 CORRELATIONS BETWEEN VALUES

A key focus of our work is correlation between values. Namely, the question of whether choice of value i is correlated with that of value j . In humans there is a robust correlation structure where certain values are more strongly correlated than others. A standard way to represent this structure is via Multidimensional Scaling (MDS) (Borg et al., 2018), calculated as follows.

First, the matrix $C \in \mathbb{R}^{19 \times 19}$ of empirical correlation coefficients is formed. Next each of the values is embedded into \mathbb{R}^2 via MDS, such that distances in \mathbb{R}^2 best approximate the correlations. For human data, this results in an approximately circular embedding, as shown in (Schwartz and Cieciuch, 2022; Skimina et al., 2021a; Daniel and Benish-Weisman, 2019). Here we performed this analysis on the LLM data. To compare the resulting dataset to the human samples, we need to normalize for the degrees of freedom of rotation and translation. This is done via Procrustes Analysis between the human and LLM embeddings. The resulting embeddings were plotted. Then, we computed the sum of squared differences between the procrusted MDS locations of each value to the human benchmark. Larger differences indicate stronger divergence from the human samples.

4 RESULTS

The above analyses were performed for all models and prompting strategies. We checked that model responses only contained scores for the questions in the questionnaires, and that they could therefore be transformed to tabular form and analyzed. This was almost always the case except for Gemma 2 27B on the Demographic prompt at temperature zero, and we therefore do not provide result for that settings.

Value Rankings: As previously mentioned, research across samples and cultures have shown that while individual differences exist in human value priorities, there are also robust common patterns. In this section, we analyze the LLM responses and compare them to the typical ranking of human values, as discussed in Section 3.2.1.

Figure 2a shows the Spearman rank correlations between human rankings and those of the different models and prompting schemes. The results show high correlation levels (> 0.8) for many prompt-model combinations. One exception is the basic prompt with GPT, which shows very low correlation. Full rankings are provided in the Appendix for several models and prompts (Table 2, Table 3 and Table 4). These reveal that values such as Benevolence that are highly ranked in humans are indeed also highly ranked by most LLMs (e.g., ranked third and first by GPT-4 for the Value Anchor prompt with temperatures 0 and 0.7 respectively). Conversely, values such as Power Dominance that are ranked low by humans, are ranked low by models (e.g., 19 by GPT-4 for the Value Anchor prompt). Figure 2b shows the scores corresponding to the Value Anchor prompt, when sorted according to human preferences. It can be seen that the models tend to agree with the human ordering on the low and high ranked values. Taken together, these results demonstrate that LLMs tend on average to align with the human ranking of values.

Correlations Between Values The MDS analysis (see Section 3.2.2) maps all values into \mathbb{R}^2 in a way that reflects their correlations. Here, we conduct MDS analyses for both human responses and LLM output, and then compare the results. The analyses were performed separately for each prompt, temperature, and model. Figure 3 illustrates a comparison between human MDS and Gemini Pro at temperature 0.0, for the Value Anchor and Names prompts, respectively. Notably, the disparities between Gemini Pro and GPT-4 for each prompt are minimal. GPT-4 plots and other Gemini Pro prompts are included in the Appendix as Figure 7 and Figure 8.

First, it can be seen that among humans, the values are organized in a circle in the theoretically expected order. These results were often identified over the years, and interpreted as resulting from the aspiration of individuals to maintain personal consistency in their motivations (Schwartz, 1992). Second, it can be seen that the MDS configuration resulting from the Value Anchor prompt more closely follows this order than the MDS resulting from the Names prompt. We further quantitatively compared the configurations in Table 1, by taking the mean squared difference between any pair of human and prompting method MDS matrices (i.e., matrices in $\mathbb{R}^{19 \times 2}$). It can be seen that the Value Anchor prompt demonstrated a better fit to human values than the other prompting methods.

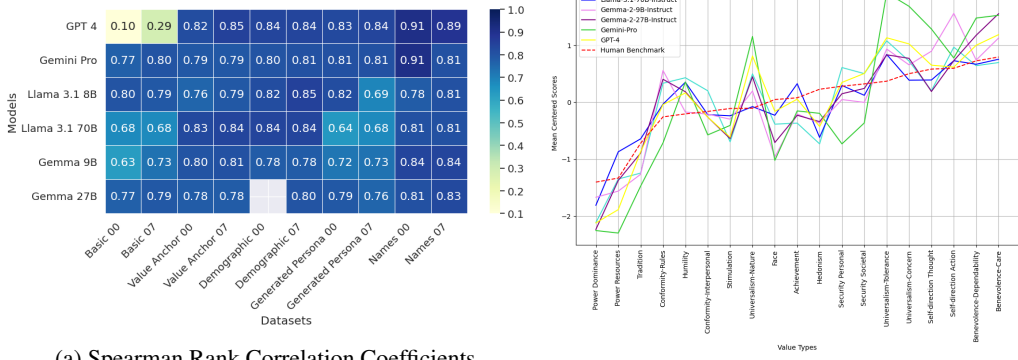


Figure 2: **Left:** A heatmap of Spearman rank correlation between benchmark value hierarchies and dataset rankings for GPT 4, Gemini Pro, Llama 3.1 8B and 70B instruct, and Gemma 2 9B and 27B across temperature conditions. **Right:** Average value scores for the Value Anchoring prompt at zero temperature. The x-axis shows values ordered according to human ranking (i.e., Power ranks lowest for humans and Benevolence ranks highest). The y-axis is the mean-centered scores the models ascribe to these values in the questionnaire, and human values in red. It can be seen that models tend to give lower scores to values that are ranked lower by humans, and higher scores to values ranked higher. The LLM scores also track the human scores (red curve) quite well.

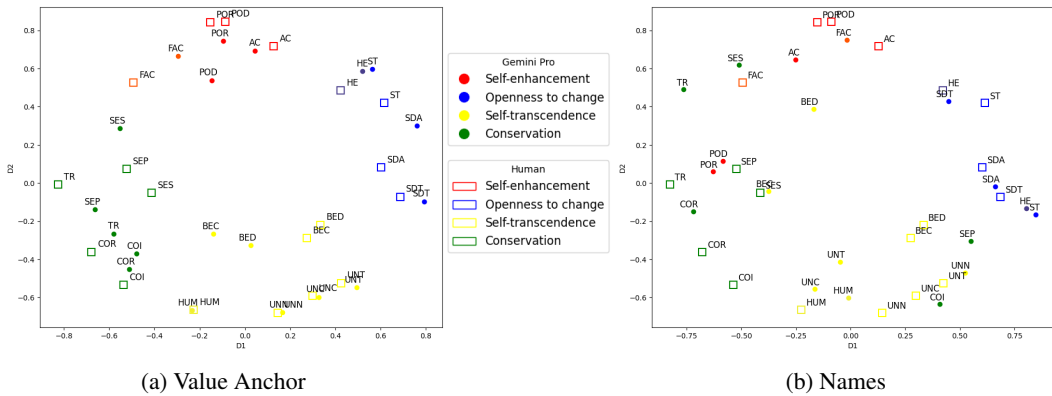


Figure 3: Comparison of Procrustes Analysis results between human data (Schwartz and Cieciuch, 2022) and Gemini Pro for Value Anchor and Names prompts, for temperature 0.0. The sum of squared differences, which measures the fit to human data, is 0.11 for the Value Anchor and 0.71 for the Names, indicating a better fit for the Value Anchor. For acronyms, refer to Section C.

4.1 UNDERSTANDING VALUE ANCHORING

The results above show that value anchoring generates a correlation structure between values that is in better agreement with that of humans. We next set out to understand why this is the case. As we shall show, anchoring on a value not only increases the score the model assigns to the value, but it only changes the way other values are scored. Specifically, values that are close to the anchored value (on the value circle, Figure 1) tend to receive high scores. On the other hand, values that are far from the anchor receive consistently lower scores. This in turn has the effect of increasing correlation between values in the anchoring setting.

To show the above, we produce an “anchored score curve” as follows. First, we order the 19 anchoring values according to their order on the value circle in Figure 1. Note that in this order, values 19 and 1 will actually be close in the circle. Then, for each set of responses of an LLM to a Value Anchor

	Basic	Value Anchor	Demographic	Persona	Names
GPT 4					
00	0.92	0.23	0.53	0.25	0.32
07	0.88	0.22	0.74	0.22	0.28
Gemini Pro					
00	0.87	0.11	0.42	0.39	0.71
07	0.69	0.11	0.75	0.28	0.57
Llama 3.1 8B					
00	0.80	0.18	0.47	0.58	0.60
07	0.57	0.16	0.47	0.58	0.57
Llama 3.1 70B					
00	0.61	0.10	0.29	0.37	0.45
07	0.44	0.10	0.22	0.40	0.44
Gemma 2 9B					
00	0.42	0.10	0.19	0.39	0.23
07	0.82	0.11	0.16	0.32	0.12
Gemma 2 27B					
00	NA	0.16	NA	0.31	0.23
07	0.64	0.17	0.15	0.25	0.19

Table 1: Sum of squared difference between the MDS embeddings of humans and LLM. Gemma 2 27B did not produce parseable results for the Demographic prompt, and for Gemma 27B at temperature 0, some values had zero-variance, thus precluding computation of correlation coefficients. All Llama models are Instruct.

prompt, we shift the anchored value to zero. We then average all these shifted curves. Results for all models are shown in Figure 4, along with a sine function which shows a good fit to these curves. The behavior of all models is quite consistent, except for Gemma-2-9B.

As expected, the anchored value receives the highest score. However, what is more interesting is that the values close to it tend to get similarly high, and values farther away (e.g., 180 degrees apart) receive the lower values. This means that the anchored model scores values in a way that is consistent with its anchoring. This in turn implies that neighboring values will tend to be correlated, thus explaining why Value Anchoring better captures human correlation patterns.

5 DISCUSSION

In the current study we analyzed the values exhibited in LLM responses. We used two metrics to estimate the quality of the LLM responses against human responses: value ranking, and value correlations. Our results highlighted the importance of the prompting mechanism. Using the Basic prompt (namely, just providing a questionnaire with no further instructions), the LLM was likely to either generate negligible variance across generated personas, or generate internally inconsistent outputs (respond differently to questions about the same value). These results suggest that LLMs cannot be treated as ‘individuals’ holding a coherent set of value priorities. In contrast, our results indicated high consistency across model types, including commercial and open LLMs.

Prompts that endow the LLM with a “personality” improved the consistency of each specific value profile, to varying degrees. The value hierarchy was consistently found across prompts, indicating that at the mean level, LLMs can simulate value rankings of human populations. More variability between conditions was found in the measure of inter-value correlations. This is arguably the most important metric since it allows analyzing consistency across values, within a single session. We found best consistency in the Values Anchor prompt. These results suggest that LLMs, applying suitable prompts, can produce a ‘population’ of individuals, each reporting a different, but coherent set of value priorities.

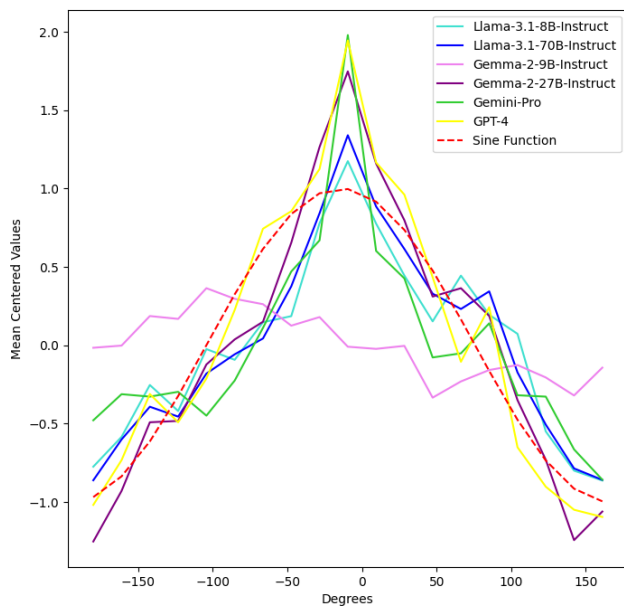


Figure 4: Analysis of scores after value anchoring. The plot shows the average of the score values after shifting to the anchored value. It can be seen that the anchored value receives the highest score, as expected. More surprisingly, neighboring values receive similarly high values, whereas more distant values receive lower values.

It is important to note that in neither of the prompts did the LLMs receive instructions for answering about all values. The results suggest that LLM is not only instruction-following, but uses the instructions as a context that guides consistent answers with relations to a variety of values. One fascinating question is where the LLM learns to produce such clear profiles of values. These profiles may be implicitly learned during pre-training. Indeed, past studies indicated that values can be identified in texts, such as newspaper articles and social media. However, these values did not necessarily follow the theoretical value inter-relations identified here (Bardi et al., 2008; Ponizovskiy et al., 2020; Kumar et al., 2018). Individuals who value competing values may experience stress and indecision when faced with a dilemma, resulting in gradual change in values toward a more coherent form (Bardi et al., 2009; Daniel and Benish-Weisman, 2019). In contrast, text may very well present both sides of a dilemma and thus retain inconsistencies. Such inconsistencies were not identified because past studies relied mostly on lexical approaches (Bardi et al., 2008; Ponizovskiy et al., 2020; Kumar et al., 2018). LLMs, taking context into account, may be more likely to identify value inter-relations correctly. LLMs may also have learned to produce value profiles in the process of fine-tuning or RLHF (Qiu et al., 2022). Future research can try to distinguish these two sources of learning using careful analysis of training sources as well as evaluation of different checkpoints in the training process.

Past research into human personas sought ways to estimate the ability of the LLM to maintain a consistent persona across a conversation (Wang et al., 2024). We establish that the unique qualities of human values, and the ample empirical knowledge collected about them, allow their use as such a method (Sagiv and Schwartz, 2022; Sagiv et al., 2017; Knafo-Noam et al., 2024). We suggest that known behavioral correlations in humans can be applied to assess the consistency of LLM personas. Here we focused on evaluations via a questionnaire, but one could envision more elaborate evaluations that rely on other features of human personalities.

The procedures and data produced here may have important contributions for psychological research. Investigators interested in human behavior can apply these procedures to produce datasets that simulate human samples. Future research can investigate their possible use to replicate known findings (e.g., age differences in values) or pretest novel hypotheses (e.g. associations between values and specific behaviors). The use of both commercial and open LLMs increases the reproducibility of the results.

The question of values in LLMs is of course of philosophical and societal importance. Our results show that on average, these values largely reproduce international value rankings. However, small variations in value importance may have implications at the societal and individual level (e.g. gender roles (Lomazzi and Seddig, 2020); entrepreneurship (Woodside et al., 2020); prosocial behavior (Daniel et al., 2020); and antisocial behavior (Benish-Weisman, 2019)). Future work should consider the influence of these values on LLM responses, as well as on the individuals interacting with them.

The current study focused on a limited number of contexts (e.g., five prompts, two temperatures, and six models). Importantly, we found commonalities across the various contexts, beside the differences. Future studies can use these results to understand what other contexts should be investigated, to possibly further enhance the quality of the output. Another limitation is the restriction to one questionnaire (PVQ-RR) corresponding a specific value system. It will be interesting to explore other forms of probing values.

REFERENCES

- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. Using large language models to simulate multiple humans and replicate human subject studies. In *International Conference on Machine Learning*, pages 337–371. Proceedings of Machine Learning Research, 2023.
- Ronald Fischer, Markus Luczak-Roesch, and Johannes A Karl. What does chatgpt return about human values? exploring value bias in chatgpt using a descriptive value theory. *arXiv preprint*, 2023.
- Lilach Sagiv and Shalom H Schwartz. Personal values across cultures. *Annual review of psychology*, 73(1):517–546, 2022. doi: 10.1146/annurev-psych-020821-125100.
- Lilach Sagiv, Sonia Roccas, Jan Cieciuch, and Shalom H Schwartz. Personal values in human life. *Nature human behaviour*, 1(9):630–639, 2017. doi: 10.1038/s41562-017-0185-3.
- Brent W Roberts and Hee J Yoon. Personality psychology. *Annual Review of Psychology*, 73(1): 489–516, 2022. doi: 10.1146/annurev-psych-020821-114927.
- Shalom H Schwartz. 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, 1992.
- Shalom H Schwartz. An overview of the Schwartz theory of basic values. *Online readings in Psychology and Culture*, 2(1):1–20, 2012. doi: 10.9707/2307-0919.1116.
- Shalom H Schwartz and Anat Bardi. Value hierarchies across cultures: Taking a similarities perspective. *Journal of cross-cultural Psychology*, 32(3):268–290, 2001. doi: 10.1177/0022022101032003002.
- Ewa Skimina, Jan Cieciuch, and William Revelle. Between-and within-person structures of value traits and value states: Four different structures, four different interpretations. *Journal of Personality*, 89(5):951–969, 2021a.
- Ella Daniel, Anna K Döring, and Jan Cieciuch. Development of intraindividual value structures in middle childhood: A multicultural and longitudinal investigation. *Journal of Personality*, 91(2): 482–496, 2023.
- Shalom H Schwartz and Jan Cieciuch. Measuring the refined theory of individual values in 49 cultural groups: psychometrics of the revised portrait value questionnaire. *Assessment*, 29(5):1005–1019, 2022. doi: 10.1177/1073191121998760.
- Shalom H Schwartz, Jan Cieciuch, Michele Vecchione, Eldad Davidov, Ronald Fischer, Constanze Beierlein, Alice Ramos, Markku Verkasalo, Jan-Erik Lönnqvist, Kursad Demirutku, et al. Refining the theory of basic individual values. *Journal of personality and social psychology*, 103(4): 663–688, 2012. doi: 10.1037/a0029393.
- Shalom H Schwartz. The refined theory of basic values. *Values and behavior: Taking a cross cultural perspective*, pages 51–72, 2017.

-
- Ali Pakizeh, Jochen E Gebauer, and Gregory R Maio. Basic human values: Inter-value structure in memory. *Journal of Experimental Social Psychology*, 43(3):458–465, 2007. doi: 10.1016/j.jesp.2006.04.007.
- Ewa Skimina, Jan Ciecuch, and Włodzimierz Strus. Traits and values as predictors of the frequency of everyday behavior: Comparison between models and levels. *Current Psychology*, 40(1):133–153, 2021b. doi: 10.1007/s12144-018-9892-9.
- Caroline Lindahl and Helin Saeid. Unveiling the values of ChatGPT: An explorative study on human values in AI systems, 2023.
- Marilù Miotto, Nicola Rossberg, and Bennett Kleinberg. Who is GPT-3? an exploration of personality, values and demographics. *arXiv*, 2022.
- Nino Scherrer, Claudia Shi, Amir Feder, and David M. Blei. Evaluating the moral beliefs encoded in llms. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Dorit Hadar-Shoval, Kfir Asraf, Yonathan Mizrahi, Yuval Haber, and Zohar Elyoseph. Assessing the alignment of large language models with human values for mental health integration: Cross-sectional study using schwartz’s theory of basic values. *JMIR Mental Health*, 11, 2024.
- Grgur Kovač, Masataka Sawayama, Rémy Portelas, Cédric Colas, Peter Ford Dominey, and Pierre-Yves Oudeyer. Large language models as superpositions of cultural perspectives. *arXiv preprint arXiv:2307.07870*, 2023.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023. doi: 10.1145/3560815.
- Dorit Hadar-Shoval, Zohar Elyoseph, and Maya Lvovsky. The plasticity of chatgpt’s mentalizing abilities: Personalization for personality structures. *Frontiers in Psychiatry*, 14:1234397, 2023. doi: 10.3389/fpsy.2023.1234397.
- Hang Jiang, Xiajie Zhang, Xubo Cao, Jad Kabbara, and Deb Roy. Personallm: Investigating the ability of gpt-3.5 to express personality traits and gender differences. *arXiv*, 2023.
- Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. In-context impersonation reveals large language models’ strengths and biases. *Advances in Neural Information Processing Systems*, 36, 2024.
- Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and David Wingate. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, 2023.
- Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*, 2023.
- Hua Li, Yu Quan Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Michael Lewis, and Katia Sycara. Theory of mind for multi-agent collaboration via large language models. *arXiv preprint arXiv:2310.10701*, 2023.
- Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. Supervised contrastive learning for pre-trained language model fine-tuning. *arXiv preprint arXiv:2011.01403*, 2020.
- Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in chatgpt. *Nature Computational Science*, 3(10):833–838, 2023. doi: 10.1038/s43588-023-00527-x.
- Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6), 2023. doi: 10.1073/pnas.2218523120.

-
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Claire Stevenson, Iris Smal, Matthijs Baas, Raoul Grasman, and Han van der Maas. Putting gpt-3’s creativity to the (alternative uses) test. *arXiv preprint arXiv:2206.08932*, 2022.
- Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. Toxicity in chatgpt: Analyzing persona-assigned language models. *arXiv preprint*, 2023.
- Xintao Wang, Yunze Xiao, Jen tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, Jiangjie Chen, Cheng Li, and Yanghua Xiao. Incharacter: Evaluating personality fidelity in role-playing agents through psychological interviews. *arXiv preprint arXiv:2310.17976*, 2024.
- Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. Bias runs deep: Implicit reasoning biases in persona-assigned llms. *arXiv preprint arXiv:2311.04892*, 2024.
- Julie A Lee, Joanne N Sneddon, Timothy M Daly, Shalom H Schwartz, Geoffrey N Soutar, and Jordan J Louviere. Testing and extending schwartz refined value theory using a best–worst scaling approach. *Assessment*, 26(2):166–180, 2019. doi: 10.1177/1073191116683799.
- National Academies of Sciences, Engineering, and Medicine. *Measuring Sex, Gender Identity, and Sexual Orientation*. National Academies Press, Washington, DC, 2022.
- Christian Haerpfer, Ronald Inglehart, Alejandro Moreno, Christian Welzel, Kseniya Kizilova, Jaime Diez-Medrano, Marta Lagos, Pippa Norris, Eduard Ponarin, and Bi Puranen, editors. *World Values Survey: Round Seven - Country-Pooled Datafile Version 5.0*. JD Systems Institute & WVSA Secretariat, Madrid, Spain & Vienna, Austria, 2022. doi: 10.14281/18241.20.
- Alison Laver-Fawcett, Leanne Brain, Courtney Brodie, Lauren Cardy, and Lisa Manaton. The face validity and clinical utility of the activity card sort–united kingdom (acs-uk). *British Journal of Occupational Therapy*, 79(8):492–504, 2016. doi: 10.1177/0308022616629167.
- Myra Cheng, Esin Durmus, and Dan Jurafsky. Marked personas: Using natural language prompts to measure stereotypes in language models. *arXiv preprint arXiv:2305.18189*, 2023.
- Ingwer Borg, Patrick JF Groenen, and Patrick Mair. *Applied multidimensional scaling and unfolding*. Springer Science & Business Media, New York, NY, 2nd edition, 2018. doi: <https://doi.org/10.1007/978-3-319-73471-2>.
- Ella Daniel and Maya Benish-Weisman. Value development during adolescence: Dimensions of change and stability. *Journal of personality*, 87(3):620–632, 2019.
- Anat Bardi, Rachel M Calogero, and Brian Mullen. A new archival approach to the study of values and value–behavior relations: validation of the value lexicon. *Journal of Applied Psychology*, 93(3):483, 2008.
- Vladimir Ponizovskiy, Murat Ardag, Lusine Grigoryan, Ryan Boyd, Henrik Dobewall, and Peter Holtz. Development and validation of the personal values dictionary: A theory–driven tool for investigating references to basic human values in text. *European Journal of Personality*, 34(5): 885–902, 2020.
- Upendra Kumar, Aishwarya N Reganti, Tushar Maheshwari, Tanmoy Chakroborty, Björn Gambäck, and Amitava Das. Inducing personalities and values from language use in social network communities. *Information Systems Frontiers*, 20:1219–1240, 2018.
- Anat Bardi, Julie Anne Lee, Nadi Hofmann-Towfigh, and Geoffrey Soutar. The structure of intraindividual value change. *Journal of personality and social psychology*, 97(5):913, 2009.

-
- Liang Qiu, Yizhou Zhao, Yuan Liang, Pan Lu, Weiyan Shi, Zhou Yu, and Song-Chun Zhu. Towards socially intelligent agents with mental state transition and human value. In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 146–158, 2022.
- Ariel Knafo-Noam, Ella Daniel, and Maya Benish-Weisman. The development of values in middle childhood: Five maturation criteria. *Current Directions in Psychological Science*, 33(1):18–26, 2024.
- Vera Lomazzi and Daniel Seddig. Gender role attitudes in the international social survey programme: Cross-national comparability and relationships to cultural values. *Cross-Cultural Research*, 54(4): 398–431, 2020.
- Arch G Woodside, Carol M Megehee, Lars Isaksson, and Graham Ferguson. Consequences of national cultures and motivations on entrepreneurship, innovation, ethical behavior, and quality-of-life. *Journal of Business & Industrial Marketing*, 35(1):40–60, 2020.
- Ella Daniel, Maya Benish-Weisman, Joanne N Sneddon, and Julie A Lee. Value profiles during middle childhood: Developmental processes and social behavior. *Child Development*, 91(5): 1615–1630, 2020.
- Maya Benish-Weisman. What can we learn about aggression from what adolescents consider important in life? the contribution of values theory to aggression research. *Child Development Perspectives*, 13(4):260–266, 2019.

Table 2: Comparative analysis of 19 values’ relative importance of the Value Anchor and Names datasets across temperatures for GPT-4 and Gemini Pro.

Benchmark		GPT-4												Gemini Pro													
		Human Data			Value Anchor 00			Value Anchor 07			Names 00			Names 07			Value Anchor 00			Value Anchor 07			Names 00			Names 07	
Rank	Values	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank
1	BEC	0.79	3	1.18	1	1.24	1	0.66	3	0.65	3.5	1.52	3	1.44	3	1.56	3	-0.10	13								
2	BED	0.72	5	0.92	4	1.00	4	0.66	3	0.65	6	1.39	4	1.30	4	1.48	5	0.03	12								
3	SDA	0.60	7	0.59	7	0.55	7	0.66	3	0.65	3.5	0.88	7	0.79	7	1.20	7	0.93	1								
4	SDT	0.58	6	0.70	6	0.69	6	0.66	3	0.65	1.5	1.33	5	1.17	5	1.48	4	0.47	5								
5	UNC	0.50	3	1.12	3	1.07	3	0.66	3	0.65	1.5	1.68	2	1.63	2	1.68	1	0.74	3								
6	UNT	0.37	1	1.26	2	1.20	2	0.66	6	0.65	5	1.93	1	1.86	1	1.63	2	0.45	6								
7	SES	0.32	8	0.41	8	0.36	8	0.64	7	0.58	7	-0.49	12	-0.42	12	0.60	8	0.14	10								
8	SEP	0.28	9	0.21	9	0.26	9	0.48	11	0.43	9	-0.87	14	-0.75	14	-0.75	14	0.21	8								
9	HE	0.23	14	-0.31	14	-0.33	14	0.50	9	0.41	10	0.04	9	-0.02	9	0.04	10	0.57	4								
10	AC	0.08	11	0.07	11	0.09	11	0.06	14	0.10	14	-0.08	11	-0.22	11	-0.35	11	-0.54	16								
11	FAC	0.05	13	-0.28	13	-0.31	13	0.30	12	0.20	13	-1.01	16	-0.91	16	-0.68	13	0.10	11								
12	UNN	-0.10	4	0.98	4	0.97	4	0.58	8	0.51	8	1.17	6	1.10	6	1.29	6	0.19	9								
13	ST	-0.11	16	-0.48	16	-0.44	16	-0.10	15	-0.14	15	-0.07	10	-0.09	10	-0.68	12	-1.03	18								
14	COI	-0.16	15	-0.41	15	-0.42	15	-0.74	17	-0.65	17	-0.72	13	-0.58	13	-0.90	16	0.34	7								
15	HUM	-0.20	10	0.20	10	0.20	10	0.30	13	0.22	12	0.12	8	0.15	8	0.52	9	0.84	2								
16	COR	-0.26	12	-0.21	12	-0.18	12	0.48	10	0.41	11	-0.95	15	-0.79	15	-1.18	17	-0.33	14								
17	TR	-0.72	17	-0.98	17	-0.93	17	-0.62	16	-0.59	16	-1.57	17	-1.43	17	-0.88	15	-0.52	15								
18	POR	-1.33	18	-1.91	18	-1.83	18	-2.60	18	-2.58	18	-2.19	19	-2.12	19	-3.09	19	-0.83	17								
19	POD	-1.40	19	-2.14	19	-2.05	19	-2.82	19	-2.81	19	-2.14	18	-2.10	18	-2.98	18	-1.64	19								

Table 3: Comparative analysis of values' relative importance for the Llama 3.1 8B and Llama 3.1 70B datasets across temperatures.

Benchmark		Llama 3.1 8B												Llama 3.1 70B															
		Human Data			Value Anchor 00			Value Anchor 07			Names 00			Names 07			Value Anchor 00			Value Anchor 07			Names 00			Names 07			
		Rank	Values	Mean	Mean	Rank	Rank	Mean	Mean	Rank	Rank	Mean	Mean	Rank	Rank	Mean	Mean	Rank	Rank	Mean	Mean	Rank	Rank	Mean	Mean	Rank	Rank	Mean	Mean
1	BEC	0.79	0.74	4	0.87	3	0.62	8	0.82	4	0.80	2	0.80	2	0.87	3	0.83	5											
2	BED	0.72	0.64	5	0.64	5	0.62	9	0.65	7	0.64	4	0.69	4	0.67	7	0.69	7											
3	SDA	0.60	1.00	2	1.04	2	1.12	2	1.12	2	0.76	3	0.72	3	0.84	5	0.84	4											
4	SDT	0.58	0.24	11	0.34	9	0.89	5	0.67	6	0.48	5	0.53	5	1.01	2	0.96	2											
5	UNC	0.50	0.77	3	0.80	4	1.11	3	0.99	3	0.46	6	0.49	6	0.85	4	0.84	3											
6	UNT	0.37	1.14	1	1.12	1	1.19	1	1.20	1	0.90	1	0.88	1	1.10	1	1.10	1											
7	SES	0.32	0.34	9	0.47	8	0.91	4	0.81	5	-0.01	11	0.01	10	0.05	11	0.08	11											
8	SEP	0.28	0.62	6	0.58	6	0.84	6	0.64	8	0.23	9	0.23	9	0.34	8	0.39	8											
9	HE	0.23	-0.61	16	-0.58	16	-1.38	16	-0.50	14	-0.51	16	-0.46	16	-0.13	13	-0.31	13											
10	AC	0.08	-0.37	13	-0.36	13	-0.27	13	-1.20	12	0.32	8	0.26	8	0.23	10	0.24	9											
11	FAC	0.05	-0.46	14	-0.58	15	0.75	15	-0.63	16	-0.24	14	-0.35	14	-0.60	15	-0.66	15											
12	UNN	-0.10	0.55	7	0.53	7	0.56	11	0.48	9	0.01	10	-0.03	11	0.29	9	0.15	10											
13	ST	-0.11	-0.59	15	-0.53	14	-0.67	14	-0.34	13	-0.15	13	-0.12	13	-0.03	12	-0.08	12											
14	COI	-0.16	-0.00	12	-0.14	12	0.61	10	0.29	10	-0.30	15	-0.37	15	-0.89	17	-0.74	16											
15	HUM	-0.20	0.44	8	0.29	10	0.69	7	0.47	10	0.35	7	0.29	7	0.74	6	0.81	6											
16	COR	-0.26	0.29	10	0.22	11	0.56	12	0.03	11	-0.14	12	-0.11	12	-0.35	14	-0.36	14											
17	TR	-0.72	-1.28	17	-1.28	17	-2.01	18	-1.53	17	-0.68	17	-0.66	17	-0.82	16	-0.82	17											
18	POR	-1.33	-1.35	18	-1.29	18	-1.77	17	-1.55	18	-0.89	18	-0.88	18	-1.46	18	-1.46	18											
19	POD	-1.40	-2.12	19	-2.11	19	-2.88	19	-2.56	19	-1.86	19	-1.81	19	-2.39	19	-2.36	19											

Table 4: Comparative analysis of values' relative importance for the Gemma 2 9B and Gemma 2 27B datasets across temperatures.

Benchmark		Gemma 2 9B												Gemma 2 27B														
		Human Data			Value Anchor 00			Value Anchor 07			Names 00			Names 07			Value Anchor 00			Value Anchor 07			Names 00			Names 07		
		Rank	Values	Mean	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank	Mean	Rank
1	BEC	0.79	1.14	2	1.22	2	1.90	2	1.88	2	1.57	1	1.57	1	1.53	1	1.50	1	1.50	1	1.50	1	1.50	1	1.50	1	1.50	1
2	BED	0.72	0.73	5	0.78	5	1.29	4	1.24	4	1.16	2	1.17	2	1.29	3	1.25	3	1.25	3	1.25	3	1.25	3	1.25	3	1.25	3
3	SDA	0.60	1.55	1	1.58	1	1.77	3	1.77	3	0.84	4	0.84	4	1.50	2	1.46	2	1.46	2	1.46	2	1.46	2	1.46	2	1.46	2
4	SDT	0.58	0.90	4	0.89	3	1.93	1	1.91	1	0.30	8	0.27	8	1.20	4	1.18	4	1.18	4	1.18	4	1.18	4	1.18	4	1.18	4
5	UNC	0.50	0.67	6	0.65	6	1.03	6	1.05	5	0.83	5	0.81	5	1.02	6	1.04	5	1.04	5	1.04	5	1.04	5	1.04	5	1.04	5
6	UNT	0.37	0.94	3	0.84	4	1.03	5	1.02	6	0.92	3	0.89	3	1.02	5	1.02	6	1.02	6	1.02	6	1.02	6	1.02	6	1.02	6
7	SES	0.32	-0.01	10	0.00	10	0.75	8	0.76	8	0.11	10	0.12	10	0.32	10	0.33	9	0.33	9	0.33	9	0.33	9	0.33	9	0.33	9
8	SEP	0.28	0.08	9	0.05	9	0.15	9	0.12	9	0.07	11	0.07	11	0.11	11	0.07	11	0.07	11	0.07	11	0.07	11	0.07	11	0.07	11
9	HE	0.23	-0.33	15	-0.32	15	-0.54	13	-0.55	13	-0.23	13	-0.23	13	-0.55	13	-0.49	13	-0.49	13	-0.49	13	-0.49	13	-0.49	13	-0.49	13
10	AC	0.08	-0.22	13	-0.22	13	-0.50	12	-0.50	12	-0.19	12	-0.21	12	-0.34	12	-0.32	12	-0.32	12	-0.32	12	-0.32	12	-0.32	12	-0.32	12
11	FAC	0.05	-1.01	16	-1.05	16	-1.45	17	-1.45	17	-0.84	16	-0.81	16	-1.24	17	-1.24	17	-1.24	17	-1.24	17	-1.24	17	-1.24	17	-1.24	17
12	UNN	-0.10	0.21	8	0.20	8	-0.03	10	0.01	10	0.56	6	0.45	6	0.55	7	0.61	7	0.61	7	0.61	7	0.61	7	0.61	7	0.61	7
13	ST	-0.11	-0.27	14	-0.24	14	-0.72	14	-0.77	14	-0.49	15	-0.42	15	-0.72	14	-0.77	15	-0.77	15	-0.77	15	-0.77	15	-0.77	15	-0.77	15
14	COI	-0.16	-0.21	12	-0.20	12	-0.94	15	-0.93	15	-0.33	14	-0.31	14	-1.00	16	-1.01	16	-1.01	16	-1.01	16	-1.01	16	-1.01	16	-1.01	16
15	HUM	-0.20	-0.14	11	-0.18	11	-0.40	11	-0.41	11	0.13	9	0.13	9	0.47	8	0.47	8	0.47	8	0.47	8	0.47	8	0.47	8	0.47	8
16	COR	-0.26	0.55	7	0.63	7	0.84	7	0.85	7	0.30	7	0.29	7	0.34	9	0.31	10	0.31	10	0.31	10	0.31	10	0.31	10	0.31	10
17	TR	-0.72	-1.23	17	-1.24	17	-0.99	16	-0.96	16	-0.99	17	-0.96	17	-0.80	15	-0.76	14	-0.76	14	-0.76	14	-0.76	14	-0.76	14	-0.76	14
18	POR	-1.33	-1.63	18	-1.62	18	-2.53	18	-2.48	18	-1.41	18	-1.41	18	-2.00	18	-1.96	18	-1.96	18	-1.96	18	-1.96	18	-1.96	18	-1.96	18
19	POD	-1.40	-1.72	19	-1.77	19	-2.59	19	-2.57	19	-2.30	19	-2.26	19	-2.70	19	-2.69	19	-2.69	19	-2.69	19	-2.69	19	-2.69	19	-2.69	19

A ADDITIONAL FILES

The Python and R code used to generate our prompt sets and analyses can be made available upon request, and will be added to open-source repositories for wider public use soon.

B QUESTION VARIANTS

For each value, we use three different question variants. For example, the three question variants formulated to assess individuals' alignment with the value type "Power Dominance" (i.e., relating to asserting authority and control over others) are as follows:

- Question 6: "He desires recognition for his abilities and seeks admiration for his actions."
- Question 29: "He prefers taking charge of situations and making decisions."
- Question 41: "He actively seeks positions of power and influence, valuing control and authority over others."

C VALUE ACRONYMS

The figures in the paper use the following value acronyms: SDT = Self-Direction Thought; SDA = Self-Direction Action; ST = Stimulation; HE = Hedonism; AC = Achievement; POD = Power-Dominance; POR = Power-Resources; FAC = Face; SEP = Security-Personal; SES = Security-Societal; TR = Tradition; COR = Conformity-Rules; COI = Conformity-Interpersonal; HUM = Humility; UNN = Universalism-Nature; UNC = Universalism-Concern; UNT = Universalism-Tolerance; BEC = Benevolence-Caring; BED = Benevolence-Dependability

D EXAMPLE PORTRAIT VALUE QUESTIONNAIRE

Figure 5 provides an example for the Portrait Value Questionnaire that was used in our study.

The value measured	Item		How much like you is this person?					
	Male Version	Female Version	1	2	3	4	5	6
Self-direction Thought	It is important to him to form his views independently.	It is important to her to form her views independently.					✓	
Security Societal	It is important to him to have a strong state that can defend its citizens.	It is important to her to have a strong state that can defend its citizens.			✓			
Hedonism	It is important to him to have a good time.	It is important to her to have a good time.				✓		
Conformity-Interpersonal	It is important to him never to annoy anyone.	It is important to her never to annoy anyone.	✓					

Figure 5: Portrait Value Questionnaire—Revised - example items. The instructions provided were: "Here we briefly describe some people. Please read each description and think about how much each person is or is not like you. Tick the box to the right that shows how much the person in the description is like you". Rankings correspond to the following descriptions: 1-Not like me at all, 2-Not like me, 3-A little like me, 4-Somewhat like me, 5-Like me, 6-Very much like me.

E THE COMPLETE ITEM LIST OF BEST-WORST REFINED VALUES (BWVR)

In our value anchoring approach, we used the description of values in [Lee et al. \(2019\)](#) to prompt the LLMs. The set of descriptions is provided below.

1. **Self-direction-thought**: developing your own original ideas and opinions
2. **Self-direction-action**: being free to act independently
3. **Stimulation**: having an exciting life; having all sorts of new experiences
4. **Hedonism**: taking advantage of every opportunity to enjoy life’s pleasures
5. **Achievement**: being ambitious and successful
6. **Power-dominance**: having the power that money and possessions can bring
7. **Power-resources**: having the authority to get others to do what you want
8. **Face**: protecting your public image and avoiding being shamed
9. **Security-personal**: living and acting in ways that ensure that you are personally safe and secure
10. **Security-societal**: living in a safe and stable society
11. **Tradition**: following cultural family or religious practices
12. **Conformity-rules**: obeying all rules and laws
13. **Conformity-interpersonal**: making sure you never upset or annoy others
14. **Humility**: being humble and avoiding public recognition
15. **Benevolence-dependability**: being a completely dependable and trustworthy friend and family member
16. **Benevolence-caring**: helping and caring for the wellbeing of those who are close
17. **Universalism-concern**: caring and seeking justice for everyone especially the weak and vulnerable in society
18. **Universalism-nature**: protecting the natural environment from destruction or pollution
19. **Universalism-tolerance**: being open-minded and accepting of people and ideas, even when you disagree with them
20. **Animal welfare**: caring for the welfare of animals

F ADDITIONAL RESULTS FOR VALUE RANKINGS

Table 2 provides additional results on value rankings for several prompting approaches.

G ADDITIONAL MDS PLOTS

In the main text we provided the MDS plots for Gemini Pro for Value Anchor and Names. Here we provide further plots for Gemini Pro in Figure 7, all the GPT 4 plots in Figure 8, all of the Llama 3.1 8B plots in Figure 9 and all of Llama 3.1 70B plots Figure 10, and all of Gemma 2 9B plots Figure 11 and Gemma 2 27B plots Figure 12 for temperature 0.0.

H COMPARING BATCH AND SEQUENTIAL PROMPTING

In the main text, we focused exclusively on batch prompting, where all items from the questionnaire were presented in a single prompt. An alternative is to present the questions in sequence, and ask the model to answer a question as soon as it is presented. To investigate potential differences between batch and sequential prompting, we evaluated on Llama models (in commercial models, sequential prompting is more expensive than batch). The value-ranking results are summarized in Figure 6, and the value-correlation results in Table 5. Regarding the values rankings, significant differences were observed between two datasets in some instances (e.g., Llama 3.1 70B, Basic: $z = -2.64$, $p = .004$), while non-significant results were noted in other cases (e.g., Llama 3.1 8B, Value Anchor: $z =$

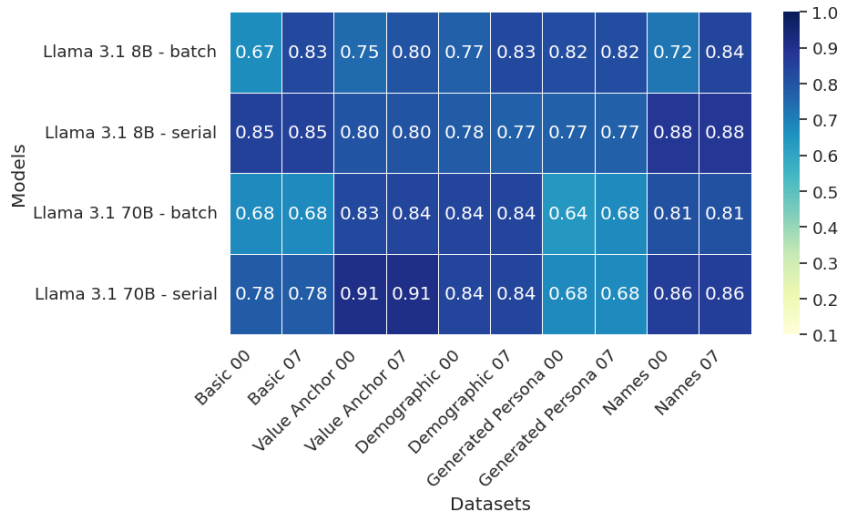


Figure 6: A heatmap of Spearman rank correlation between benchmark value hierarchies and dataset rankings for Llama 3.1 8B and 70B instruct for batch versus serial prompting methods, across temperature conditions.

Llama 3.1 8B				
	Value Anchor	Demographic	Generated Persona	Names
Batch prompting				
00	0.18	0.47	0.58	0.60
07	0.16	0.47	0.58	0.57
Serial prompting				
00	0.18	0.54	0.65	0.61
07	0.18	0.54	0.65	0.37
Llama 3.1 70B				
	Value Anchor	Demographic	Generated Persona	Names
Batch prompting				
00	0.10	0.29	0.37	0.45
07	0.10	0.22	0.40	0.44
Serial prompting				
00	0.14	0.26	0.20	0.48
07	0.14	0.26	0.20	0.69

Table 5: Sum of squared difference for MDS embeddings of humans and LLM.

-0.83, $p = .203$). This indicates that, overall, the ranking correlations are closely aligned, with no clear inclination toward either batch or sequential prompting as better replicating the human value hierarchy. As for the value-correlations, it can be seen that the sequential prompts replicate the finding that the Value Anchor prompt best captures the circular structure of human values. Interestingly, for Llama 3.1 8B, batch prompting appeared to yield superior results. However, for Llama 3.1 70B, this was not the case across most prompts, suggesting that batch prompting may not consistently perform better across different models.

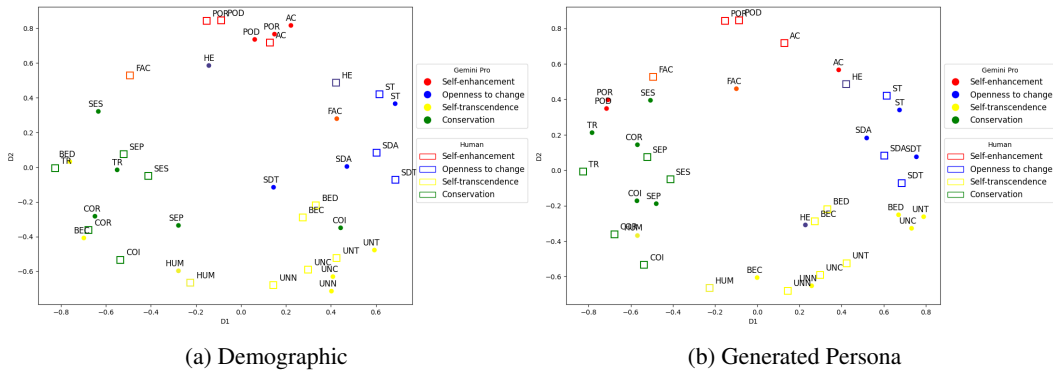


Figure 7: Comparison of the MDS results between human data (Schwartz and Cieciuch, 2022) and Gemini Pro for Demographic and Generated Persona respectively, in the temperature 0.0 condition.

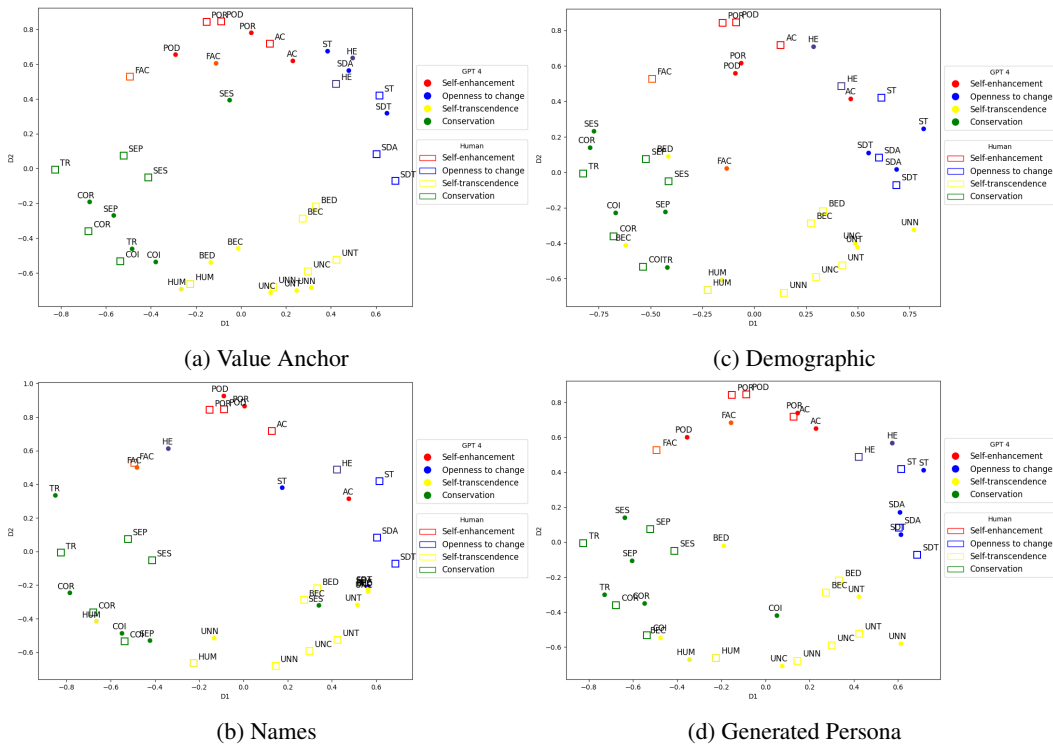


Figure 8: Comparison of the MDS results between human data (Schwartz and Cieciuch, 2022) and GPT 4 for all prompts, in the temperature 0.0 condition.

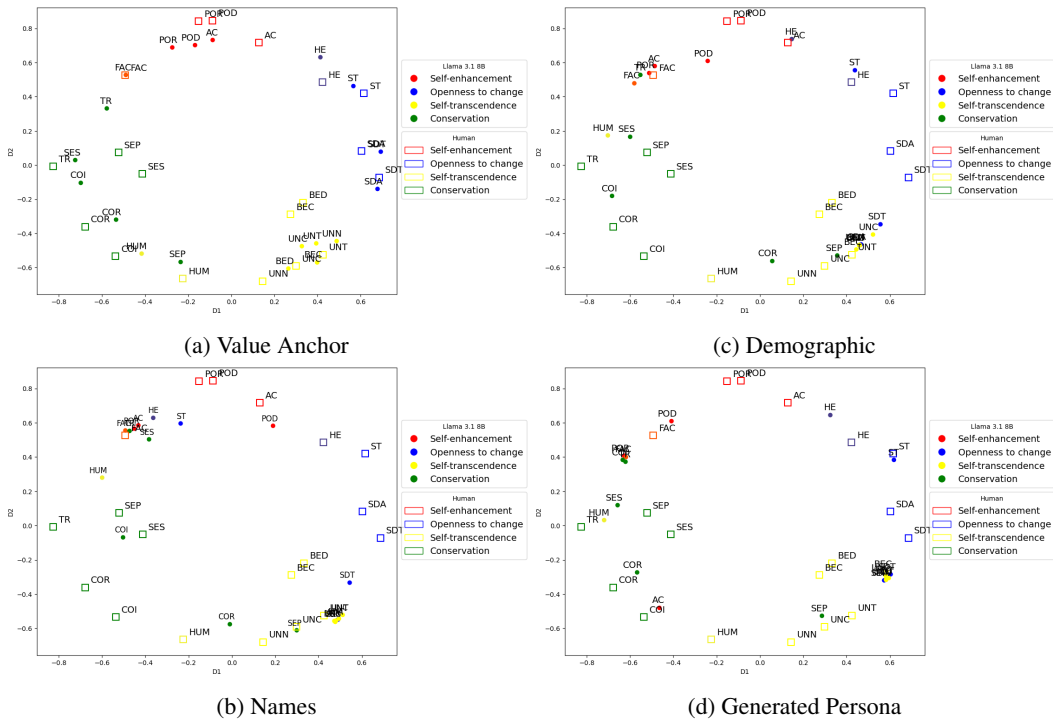


Figure 9: Comparison of the MDS results between human data (Schwartz and Ciecich, 2022) and Llama 3.1 8B for all prompts, in the temperature 0.0 condition.

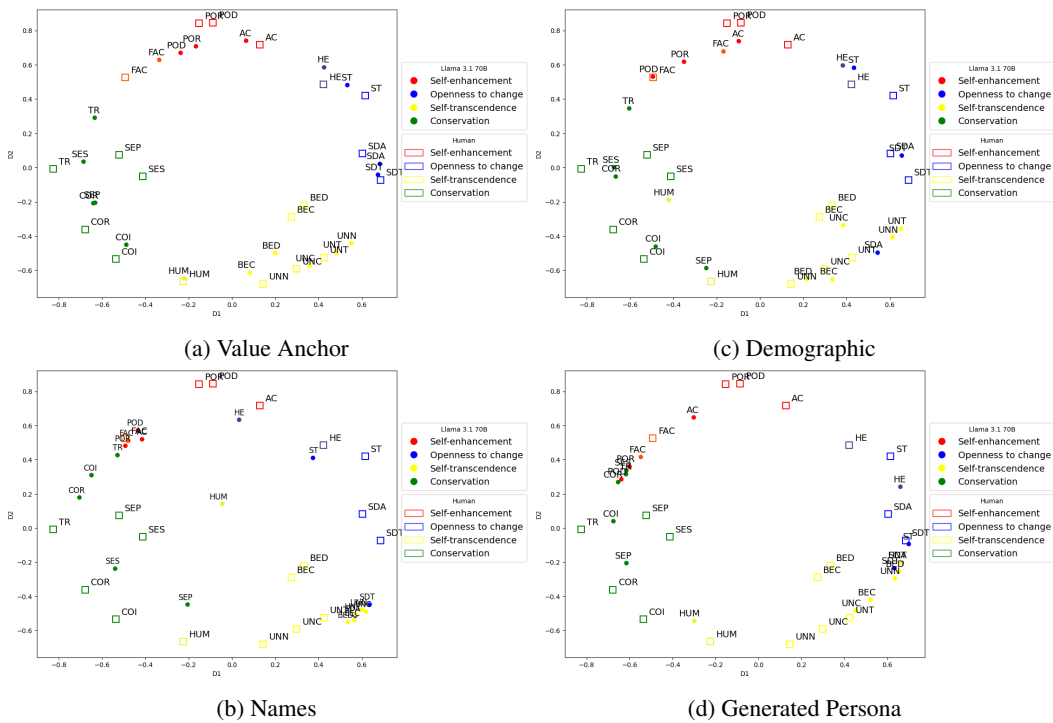


Figure 10: Comparison of the MDS results between human data (Schwartz and Ciecich, 2022) and Llama 3.1 70B for all prompts, in the temperature 0.0 condition.

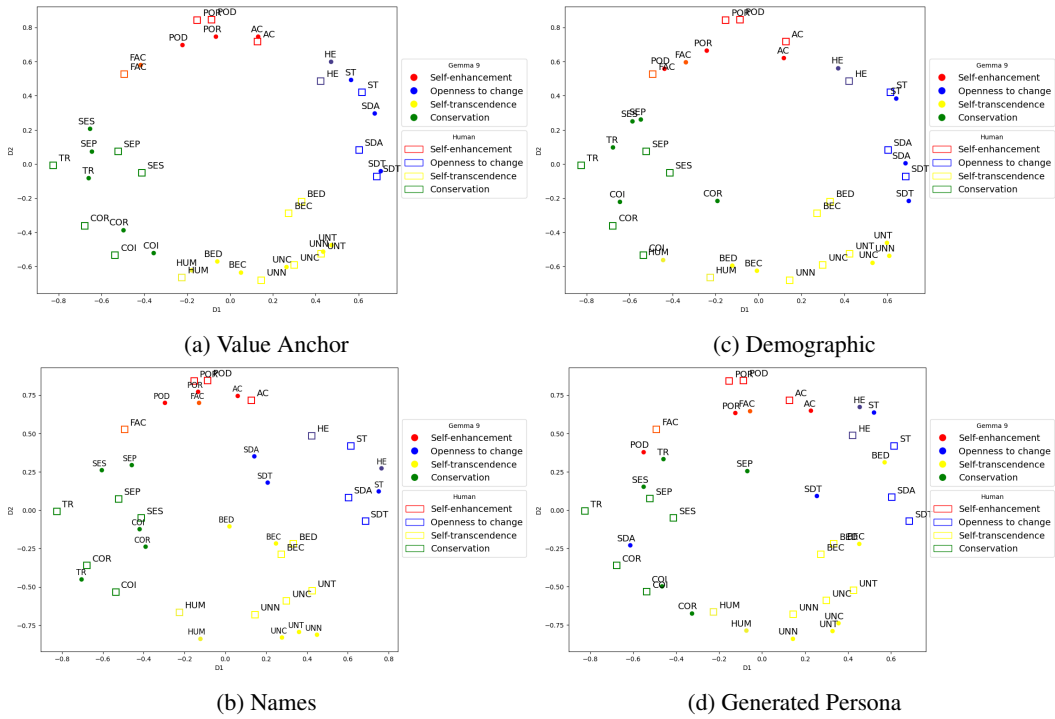


Figure 11: Comparison of the MDS results between human data (Schwartz and Cieciuch, 2022) and Gemma 2 9B for all prompts, in the temperature 0.0 condition.

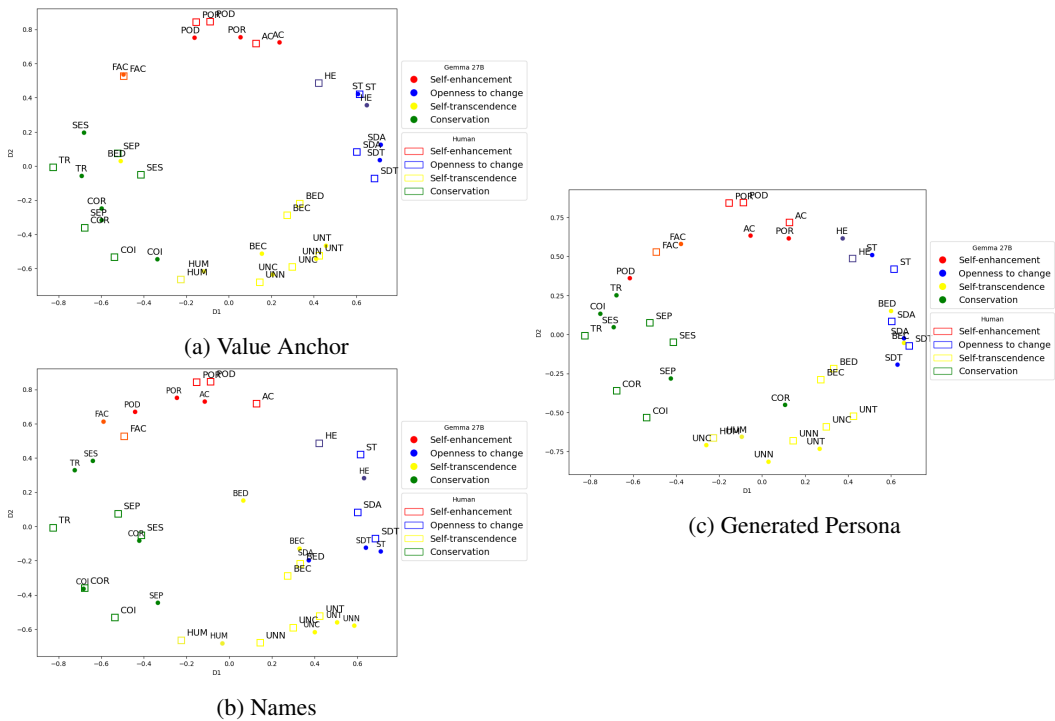


Figure 12: Comparison of the MDS results between human data (Schwartz and Cieciuch, 2022) and Gemma 2 7B for all prompts, in the temperature 0.0 condition, with the exception of the Demographic prompt-(see Footnote 2).