Investigating Cultural Alignment of Large Language Models

Badr AlKhamissi

EPFL

badr.alkhamissi@epfl.ch

Mai AlKhamissi

Anthropology, Princeton University mai.alkhamissi@princeton.edu

Abstract

The intricate relationship between language and culture has long been a subject of exploration within the realm of linguistic anthropology. Large Language Models (LLMs), promoted as repositories of collective human knowledge, raise a pivotal question: do these models genuinely encapsulate the diverse knowledge adopted by different cultures? Our study reveals that these models demonstrate greater cultural alignment along two dimensions—firstly, when prompted with the dominant language of a specific culture, and secondly, when pretrained with a refined mixture of languages employed by that culture. We quantify cultural alignment by simulating sociological surveys, comparing model responses to those of actual survey participants as references. Specifically, we replicate a survey conducted in various regions of Egypt and the United States through prompting LLMs with different pretraining data mixtures in both Arabic and English with the personas of the real respondents and the survey questions. Further analysis reveals that misalignment becomes more pronounced for underrepresented personas and for culturally sensitive topics, such as those probing social values. Finally, we introduce Anthropological Prompting, a novel method leveraging anthropological reasoning to enhance cultural alignment. Our study emphasizes the necessity for a more balanced multilingual pretraining dataset to better represent the diversity of human experience and the plurality of different cultures with many implications on the topic of cross-lingual transfer.1

1 Introduction

Large Language Models (LLMs) such as ChatGPT have garnered widespread utilization globally, engaging millions of users. Users interacting with

Muhammad ElNokrashy

Microsoft Egypt
muelnokr@microsoft.com

Mona Diab

LTI, Carnegie Mellon University mdiab@andrew.cmu.edu

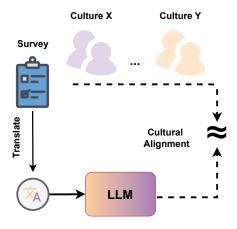


Figure 1: Our framework for measuring the cultural alignment of LLM knowledge/output and ground-truth cultural data collected through survey responses.

these models across multiple languages have observed a noteworthy phenomenon: Prompting with different languages may elicit different responses to similar queries (Lin et al., 2022; Shen et al., 2024). From our observations, one reason for the difference between the responses is that they tend to reflect the culturally specific views commonly expressed by the people which use the same language as the prompt. Here, we hypothesize that the root cause of this phenomenon lies in the training data, which encodes different and at times conflicting "knowledge" across different languages.²

Culture is a complicated term and defining it stands at the core of anthropological inquiry. Hundreds of definitions exist in literature which cover different aspects of interest (Kroeber and Kluckhohn, 1952). In this paper, we consider culture as as a multi-faceted inquiry that demonstrates substan-

Our code and data are available at https://github.com/b khmsi/cultural-trends.git

² In this work, we advocate for the term "Cultural Trends" instead of "Biases." This choice is deliberate as the term "bias" outside mathematical context often carries a negative connotation—a problematic default position. The use of Cultural Trends emphasizes that a model reflecting a particular cultural inclination does not inherently imply danger or stereotyping. Instead, it signifies alignment with the views of a specific population, highlighting cultural significance.

tial diversity among human communities, encompassing worldviews and belief systems. Through this lens, we aim to measure the cultural alignment of Large Language Models (LLMs) by simulating existing surveys that have been carried out by sociologists in specific populations. We utilize the responses from actual survey participants as our reference or gold standard. Then we measure the similarity between the model's answer when prompted with the participant's "persona" and the actual survey answer. The term "persona" in this context refers to an explicit description of a survey participant, encompassing various traits of interest such as social class, education level, and age (see Section 4.3 for a detailed description). This is done for various LLMs trained and prompted under different configurations. We use this similarity as a proxy for the degree of a model's knowledge of a particular culture. This enables us to assess the LLMs' capacity to capture the diversity not only of a specific country but also among individuals within that country.

We focus on a survey conducted in two countries: Egypt (EG) and the United States of America (US). It covers a diverse demographic set within each country with questions spanning various themes that include topics of social, cultural, material, governmental, ethical, and economic significance. This work primarily explores the impact of the language used for prompting and the language composition of pretraining data on a model's cultural alignment as defined above. We consider two languages for prompting: English and Arabic as they are the primary languages used in the surveys. Specifically, we consider four pretrained LLMs: GPT-3.5³ also known as ChatGPT, and three 13B parameter instruction-tuned models. The multilingual mT0-XXL (Muennighoff et al., 2023) is trained on a variety of languages, LLaMA-2-13B-Chat (Touvron et al., 2023) which is trained primarily on English data, and AceGPT-13B-Chat (Huang et al., 2023), a model finetuned from LLaMA-2-13B-Chat focusing on Arabic.

Our contributions include highlighting the significant role of language in the perceived, functional cultural alignment in model responses, which is affected by both (1) the language in the pretraining data and (2) that of the prompt. Further analysis shows that (3) models capture the variance of certain demographics more than others, with the gap

increasing for underrepresented groups. Finally, (4) we propose Anthropological Prompting as a method to enhance cultural alignment in LLMs.

2 Research Questions

Prompting Language and Cultural Alignment:

We hypothesize that employing the native language of a specific culture will yield greater cultural alignment compared to using a foreign language. For instance, prompting an LLM in Arabic *may* achieve higher alignment to a survey conducted in Egypt than prompting it in English.

Pretraining Data Composition: We hypothesize that, for a fixed model size, pretraining models with a higher proportion of data from a specific culture will lead to an increased alignment with the results of surveys conducted in that culture. For instance, a 13B Arabic monolingual model is expected to exhibit higher alignment than a 13B English model for a survey conducted in Egypt.

Personas and Cultural Topics: We anticipate that misalignment will increase for personas from digitally underrepresented backgrounds. For instance, alignment in both Arabic and English tests are expected to be lower for a working-class persona in Aswan (a city in the south of Egypt) compared to an upper-middle-class persona in Cairo (Egypt's capital and its most populous city). Further, we hypothesize that misalignment will increase for uncommon cultural topics.

Finetuning Models to Induce Cross-Lingual Knowledge Transfer: We gauge the effect of cross-lingual transfer for models predominantly pretrained on one language but finetuned on another. To answer this question, we use the LLaMA-2-Chat-13B model (trained primarily on an English corpus) (Touvron et al., 2023) and the AceGPT-Chat-13B model (a LLaMA-2-Chat-13B model further finetuned on a corpus of Arabic and English data) (Huang et al., 2023).

3 Anthropological Preliminaries

The concept of **culture** undergoes continual transformation, encompassing various elements that evolve with time as well as geographical and historical context. Many definitions of culture are traced back to Tylor (1871) wherein culture constitutes an integrated body of knowledge, belief, art, morals, law, custom, and any other capabilities

³ GPT-3.5 is gpt-3.5-turbo-1106 throughout this work.

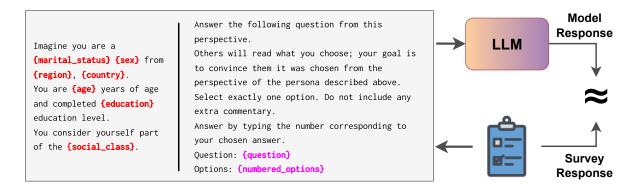


Figure 2: Template used when querying models in English. (**Left**) The model is first instructed to respond under a specific persona along the demographic parameters highlighted in red. (**Right**) The rest of the prompt instructs the model to follow the perspective of the persona closely, respond in a specific format (only the index of the answer), and avoid any extraneous commentary.

and habits expressed by members of a society. In that sense, any reflection of such aspects of life in written records can be considered a cultural trend expressed by that text. A model which expresses views in some aspect of life which is aligned with a group of people is *culturally aligned* with them in that scenario.

An alternative perspective shows culture as patterns of behavior. These patterns, how they are chosen and valued, and their meaning, manifest in different forms, such as linguistic records. In that sense, culture observes behavior through history, and affects it through population dynamics (Kroeber and Kluckhohn, 1952). The cultural expression of agential members within a society, including artificial agents such as LLMs, thus affects and is affected by the behavior and recording of ideas by fellow members. Models learn effectively from humans and equally impart their learnings upon other humans, distributing their internalized cultural ideas in the process (Clifford et al., 2020).

3.1 Working Assumptions

Given this anthropological backdrop, we describe some modeling assumptions we have adopted and the motivation behind them.

Language → Culture We assume that a language can be used as proxy for its dominant culture. Although some languages are used by multiple cultures, contemporary consideration of such languages tends to emphasize a particular culture among their diverse user base (compare the significance given to French output from France and from the Senegal). Prompting with dialects specific to a certain population can help alleviate that concern.

Culture → Language We assume that, most often, a culture will produce linguistic records and communications in one dominant language. Although, contrary to the expectation that all communications would be in the main or ostensibly official language, we know that this may not be the case. For example, individuals in Egypt may express their opinions online in English rather than in their native language for a variety of reasons.

4 Experimental Setup

4.1 World Values Survey (WVS)

The WVS project gathers responses to an array of questions on matters of social, cultural, material, governmental, ethical, and economic importance, as a rough categorization, all from demographically-controlled population samples around the world (Haerpfer et al., 2020). The latest edition (WVS-7) was conducted between 2017 and 2021. It includes some region-specific modules in addition to the globally-applied categories. WVS-7 has 259 questions and was designed to include indicators towards multiple United Nations Sustainable Development Goals. The survey is set up as a questionnaire provided to select samples from the general population. The questions in the survey are localized to the native or dominant regional languages.

In this work, we select 30 questions that encompass diverse themes. The chosen questions are intentionally not straightforward, allowing for a degree of potential cultural variation in responses. For every question, we create four linguistic variations (i.e. paraphrases) by providing ChatGPT with a short description of the question along with

Dimension	Possible Values	
Region	Cairo, Alexandria, etc.	
Sex	Male, Female	
Age	Number	
Social Class	Upper, Working, etc.	
Education Level	Higher, Middle, Lower	
Marital Status	Married, Single, etc.	

Table 1: The demographic dimensions used when prompting the model to emulate a certain survey respondent. Region is country-specific. More information in Appendix D.

the anticipated answer options from participants. The questions are translated into Arabic using machine translation, followed by manual editing by native Arabic speakers to ensure preservation of the intended meaning. More details about the generation process, including examples, are available in Appendix G.

4.2 Survey Participants

The WVS-7 survey conducted in Egypt and the United States comprised 1,200 and 2,596 participants respectively, representing diverse backgrounds. In this work, we only consider 6 demographic dimensions when prompting the LLMs. Table 1 shows the dimensions along with some possible values they can take. In addition, the left part of Figure 2 shows the template used to prompt the model in English with a specific persona. In the context of this paper, the term **persona** denotes *a singular instance of this six-dimensional tuple*.

Filtering Participants In our survey simulations, we filtered the participants to have an equal distribution across both countries along the demographic dimensions (except Region since it is country-specific). We selected participants such that for each person interviewed in Egypt we have a corresponding person who comes from exactly the same demographics from the US with the exception of the location. This resulted in 303 unique personas for each country. The distribution of the survey respondents from each country, including examples of some personas, can be found in Appendix D.

4.3 Personas: Role-Playing for LLMs

To guide a language model with instructionfollowing support in order to respond *like* a specific subject from a particular demographic, 4 we utilize personas (Joshi et al., 2023). A persona is a description of a person which covers as many traits as deemed important to be controlled for in the context of an interaction or study. Accordingly, we query the model by a prompt that specifies the values for each demographic dimension of interest. The prompt is generated from a single template and is written in ordinary prose. Figure 2 shows the template used when querying the models in English. It comprises three parts: the first specifies to the model the persona it must emulate along the 6 demographic dimensions discussed in Section 4.2. The second instructs the model to follow the perspective of the persona closely, respond in a specific format (only the index of the answer), and avoid any extraneous commentary. The last part is the question followed by a list of numbered options that the model must choose from.

4.4 Pretrained Large Language Models

Table 6 lists the models used in this work along with their corresponding number of parameters and pretraining language mixtures. In particular, we opt for instruction-tuned models as they can be assessed in a zero-shot manner by adhering to the provided instructions (Zhang et al., 2023). The largest model in our selection is GPT-3.5, primarily trained on English data; although, it has showcased competitive performance on Arabic NLP benchmarks (Alyafeai et al., 2023; Khondaker et al., 2023). The three other models are selected to be of the same size (13B parameters) for fair comparison: (1) mT0-XXL (Muennighoff et al., 2023) trained with a more balanced mixture of languages, is expected to exhibit a reduced impact of Anglocentric responses; (2) LLaMA-2-13B-Chat⁵ (Touvron et al., 2023) trained primarily on English data but is capable of responding to Arabic prompts; (3) AceGPT-13B-Chat (Huang et al., 2023) is a model finetuned on a mixture of Arabic and English data. It achieved state-of-the-art results on the Arabic Cultural and Value Alignment Dataset among opensource Arabic LLMs through localized training.

4.5 Computing Cultural Alignment

The survey simulations involve prompting each model with a specific persona, followed by an instruction and a question (refer to Figure 2). Each

⁴ A subject is a person participating in the survey.

⁵ For brevity, we omit 13B from LLaMA-2-13B-Chat and AceGPT-13B-Chat in future references.

	Egypt		Uı	nited States		
Model	English	Arabic	Ar-En	English	Arabic	En-Ar
GPT-3.5	47.08 / 23.42	50.15 / 28.56	3.07	65.95 / 40.22	63.77 / 38.36	2.18
AceGPT-Chat	46.15 / 28.83	49.49 / 30.60	3.34	54.55 / 29.94	51.12 / 25.45	3.43
LLaMA-2-Chat	47.95 / 25.61	44.67 / 23.34	-3.28	63.90 / 37.40	62.29 / 36.03	1.61
mT0-XXL	45.16 / <u>28.75</u>	<u>46.69</u> / 27.10	1.53	53.20 / 28.30	<u>57.75 / 34.51</u>	-4.55

Table 2: Cultural alignment against both Egyptian and United States survey responses using Soft / Hard similarity metrics for each model as a function of the prompting language. <u>Underlined</u> is the optimal prompting language for each model and survey. The third column in each block shows the difference in soft alignment between country's dominant language and the other language. Refer to Appendix A for results without excluding responses where equivalent personas in both surveys answered similarly.

question is independently prompted four times for each persona using the generated linguistic variations. Subsequently, we sample five responses for each question variant using a temperature of 0.7.⁶ The model's response for a particular persona and question variant is determined by computing a majority vote over the sampled responses.

Following this, we assess a model's cultural alignment by comparing its responses for each persona separately with the original subject's response in one of the two surveys. This comparison is conducted in two ways: either directly comparing the responses (Hard metric) or considering the responses while taking into account the order of the options for ordinal questions (Soft metric). We exclude instances where two subjects belonging to similar persona from both the Egypt and US surveys provided identical answers for a given question. This exclusion ensures a more accurate assessment of each model's capability in discerning the differences between the two cultures.

Hard Metric is the plain accuracy, which compares model answers to the survey responses for a given persona. Formally, the final cultural alignment is $H_{f,c} = \mathsf{mean}_{q,p} \{\mathbb{1}(\hat{y} = y)\}$, where $\hat{y}_{q,p}$ denotes the model's response after the majority vote for a question prompt q and persona p, while $y_{q,p}$ is the ground-truth response of a specific subject with persona p, all from culture c.

Soft Metric $S_{f,c}$ is a relaxed version of the hard metric which awards partial points in questions with an ordinal scale. However, if the question provides categorical options only or the subject in the survey responded with a "don't know" (regardless of the scale), the metric defaults to plain accuracy.

First we calculate the error per q, p:

$$\varepsilon_{f,c}(q,p) = \begin{cases} \frac{|\hat{y} - y|_{q,p}}{|q| - 1} & \text{IsOrd}(q,p), \\ \mathbb{1}(\hat{y} \neq y)_{q,p} & \text{otherwise} \end{cases}$$
(1)

Where $\mathsf{IsOrd}(q,p)$ refers to questions with ordinal answers where the survey subject p did not pick "don't know", and |q| is the count of options in q.

The soft alignment score for a model f is the average over all queries and personas per model and culture: $S_{f,c} = \mathsf{mean}_{q,p}\{1 - \varepsilon_{f,c}(q,p)\}$

4.6 Anthropological Prompting

Inspired by long-term ethnographic fieldwork—which stands as the primary research method within the discipline of cultural anthropology—we introduce a novel prompting method to improve cultural alignment for LLMs, **Anthropological Prompting**. The objective of engaging in extended ethnographic fieldwork is to establish meaningful connections with interlocutors, facilitating the ability to produce critical and in-depth analyses of both the subjects and the topics under study.

In this context, we strive to emulate a digital adaptation of ethnographic fieldwork by guiding the model to think as if it has been actively participating in this method. We prompt the model to comprehend the intricate complexities and nuances associated with identities, inquiries, and linguistic constructions. For instance, we elaborate on the emic and etic perspectives of examining culture, highlighting the layered nature of interpersonal connections and emphasizing how personal experiences significantly shape subjectivities. In

⁶ This was empirically set.

^{7 &}quot;Emic" refers to an insider's perspective, focusing on the internal understandings within a specific culture. Conversely, "etic" refers to an outsider's perspective.

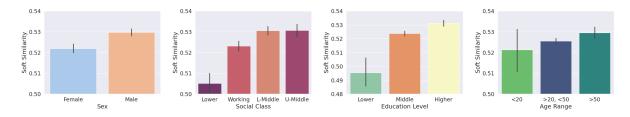


Figure 3: Cultural alignment as a function of a subject's Sex, Education Level, Social Class, and Age Range. Results are averaged across the models, prompting languages and surveys used in this work. L-Middle and U-Middle are Lower Middle and Upper Middle Class respectively.

doing so, our intention is to introduce an anthropological methodology, encouraging the model to "think" in a manner akin to an anthropologist. The exact prompt and more details about the experimental setup can be found in Appendix I.

5 Results

5.1 Anglocentric Bias in LLMs

Table 3 shows that all LLMs considered in this work—regardless of being trained to be multilingual or finetuned on culture-specific data—are significantly more culturally aligned with subjects from the US survey than those from the Egypt survey. Concurrent research has shown similar results of current LLMs exhibiting Western biases (Durmus et al., 2023; Naous et al., 2023). This can largely be attributed to the data used for training and for guiding crucial design decisions such as model architecture, tokenization scheme, evaluation methods, instruction-tuning, and so on.

Model	Egypt	United States
GPT-3.5	48.61 / 25.99	64.86 / 39.29
AceGPT-Chat	47.82 / 29.72	52.83 / 27.69
LLaMA-2-Chat	46.31 / 24.48	63.10 / 36.72
mT0-XXL	45.92 / 27.93	55.48 / 31.40
Average	47.16 / 27.03	59.07 / 33.78

Table 3: Cultural alignment against responses from both Egyptian and United States surveys using Soft / Hard similarity metrics for each model. The results are averaged across both prompting languages. The alignment with the United States populations is much higher reflecting the euro-centric bias in current LLMs.

5.2 Prompting & Pretraining Languages

Table 2 illustrates the impact of prompting language on the cultural alignment of the four LLMs examined in this study. Specifically, using each country's dominant language prompts a notable increase in alignment compared to using the alternative language for both GPT-3.5 and AceGPT-Chat, according to both metrics. For example, using Arabic to prompt both models yields better alignment with the Egypt survey than prompting with English. Conversely, English prompts result in improved alignment with the US survey compared to Arabic. However, given that LLaMA-2-Chat is predominantly pretrained on English data, we observe that Arabic prompts are less effective in enhancing alignment with the Egypt survey and thus posit that the lack of Arabic data in the pretraining leads to lack of knowledge of Egyptian culture. In contrast, for the multilingual mT0-XXL, despite being trained on a more balanced language distribution, it appears to suffer from the curse of multilinguality (Pfeiffer et al., 2022), as evidenced by its inferior cultural alignment with the US survey when prompted with English compared to Arabic. Finally, we report the models' consistency in responding to paraphrases of the same question in Appendix C.

5.3 Digitally Underrepresented Personas

Figure 3 displays the cultural alignment across various demographic variables, averaged across the four LLMs, two prompting languages, and responses from the two countries using the soft alignment metric. Surprisingly, we observe a distinct trend among the models tested in this study concerning *social class* and *education level*. Specifically, as the background of individuals changes from lower to higher levels in both respective dimensions, alignment improves. This underscores that the models better reflect the viewpoints of specific demographics over others, with marginalized populations enjoying lower alignment. Additionally, the analysis of the sex dimension reveals that the models correspond more accurately to the

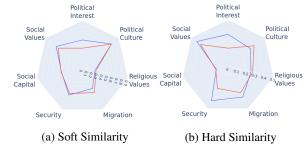


Figure 4: — Arabic — English. Alignment of GPT-3.5 with the Egypt survey using both the soft and hard metrics by theme as a function of the prompting language.

actual survey when impersonating male respondents than female respondents. Similarly, older age groups exhibit higher alignment than younger age groups.

5.4 Cultural Alignment per Theme

The 30 questions examined in this work are categorized into 7 distinct themes outlined by the WVS survey (Haerpfer et al., 2020). Table 10 illustrates the distribution of questions across these themes. The granularity provided by these themes enables us to assess alignment concerning topics such as Religious Values. In Figure 4, we illustrate the cultural alignment of GPT-3.5 with respect to responses from both the Egypt and the US survey, and examine the prompting language effect within each plot. The three themes that are contributing to the improvement in alignment in the Egypt survey when prompting in Arabic using GPT-3.5 are Social Values, Political Interest and Security. In the US survey, both English and Arabic prompting perform very closely except in the Migration theme where English has a slight edge. See Appendix H for a comprehensive set of results for all other models, metrics, and country combinations.

5.5 Finetuning for Cultural Alignment

Here, we delineate the contrast between AceGPT-Chat and LLaMA-2-Chat to illustrate the impact of finetuning an English-pretrained model on data from another language on cultural alignment. We observe an improvement in alignment with the Egypt survey across both metrics when the two models are prompted in Arabic (see Table 2 for a quantitative comparison). When prompted in English, the increase is evident only with the hard metric. Conversely, we note a decline in alignment following finetuning when evaluating alignment against the US survey, indicating that the model

Prompting Method	Soft	Hard
Vanilla	0.4834	0.2443
Anthropological	0.5102	0.2838

Table 4: Anthropological prompting outperforms Vanilla prompting across both metrics in terms of cultural alignment with the Egypt survey. Results here are on GPT-3.5 with English prompting.

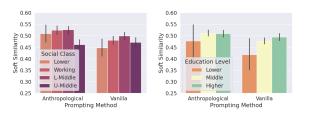


Figure 5: Anthropological prompting improves alignment for underrepresented personas compared to Vanilla prompting. Results on GPT-3.5 using English prompting. More in Appendix I.

forgot some of its existing US cultural knowledge while adapting to data in another language.

5.6 Anthropological Prompting

To improve cultural alignment with responses from Egyptian participants and underrepresented groups, we propose Anthropological Prompting. This approach enables the model to reason before answering the question while grounded with a framework adapted from the toolkit of anthropological methods. The rationale behind it is described in Section 4.6. The framework offers guidance for the model to consider emic and etic perspectives, cultural context, socioeconomic background, individual values, personal experience, cultural relativism, as well as spatial and temporal dimensions in a nuanced manner. The exact prompt is provided in Appendix I. Table 4 presents the results when prompting GPT-3.5 in English, comparing both "vanilla" (persona-based prompting) with anthropological prompting using one variant per question. While vanilla prompting generates 5 responses and computes the majority vote to determine the final answer, giving it an apparent advantage, the anthropological prompting method, which generates only one response, still outperforms it.

Further, we observe that anthropological prompting improves cultural alignment for participants from underrepresented backgrounds. Figure 5 illustrates this comparison between vanilla and anthropological prompting across Social Class and

Education Level demographic dimensions. The alignment distribution among social classes and education levels becomes more equitable as a result.

6 Discussion

In Section 5.2, we demonstrate that both the language utilized for pretraining and the language employed for prompting contribute to enhancing cultural alignment, particularly for countries where the language in question is prevalent. This observation aligns intuitively with our assumption that a culture primarily generates content in its native language on the internet.

During pretraining, a model encodes that cultural knowledge within its parameters; then during inference, the prompting language activates the subnetwork responsible for that encoded knowledge (Foroutan et al., 2022). This observation further underscores the limitation of current LLMs in effectively transferring knowledge across different languages. This is particularly evident in languages with different scripts like Arabic and English (Qi et al., 2023).

However, despite our use of Modern Standard Arabic (MSA) as the primary language for representing Egyptian culture, it is crucial to note that Egyptians do not employ MSA in their daily interactions. Hence, we posit that employing the Egyptian Arabic dialect would likely yield even greater alignment, provided that the model is sufficiently trained on this dialect. Moreover, within Egypt, there exist dialectal variations, similar to differences between various states and ethnic groups in the US. Therefore, when assessing cultural alignment, it is imperative to acknowledge the diverse backgrounds within each country; for there exists no singular, narrow Egyptian archetype, for example. This is why our study focuses on measuring personas across multiple demographic dimensions.

Finally, we would like to highlight that cultural information is often difficult to verify and coalesce into knowledge, seeing as there are numerous approaches to collecting evidence, building perspectives, and constructing theorems within cultural topics. We observe in modern large models that cultural-knowledge transfer tends to occur from a few dominant languages (for example English) and cultures into responses for prompts in other languages about other cultures. See Table 2.

7 Related Work

Measuring Subjective Opinions in LLMs: Concurrent works tackle the notion of cultural alignment but from differing perspectives. Durmus et al. (2023) similarly utilizes cross-national surveys to quantitatively assess how well LLMs capture subjective opinions from various countries. However, one notable difference from our method is that their metric solely evaluates the similarity between the model's and survey's distributions over possible options using the Jensen-Shannon Distance, without considering granularity at the persona level nor the order of options for ordinal questions.

Arora et al. (2023) measured the extent to which cross-cultural differences are encoded in multilingual encoder-only models by probing them in a cloze-style manner across multiple languages using questions from the WVS and Hofstede survey (Hofstede, 1984). Cao et al. (2023) similarly use the Hofstede Culture Survey to assess the crosscultural alignment between ChatGPT and certain societies when prompting it in different languages, showing that ChatGPT exhibits a strong alignment with American culture and adapts less effectively to other cultural contexts.

Naous et al. (2023) demonstrate that multilingual and Arabic monolingual LMs exhibit trends from Western cultures even when prompted in Arabic and contextualized within an Arabic cultural setting. Lahoti et al. (2023) propose a novel prompting method aimed at enhancing cultural diversity in LLM responses.

Tjuatja et al. (2023) demonstrate that LLMs should not be relied upon as proxies for gauging human opinions, as they do not accurately reflect response biases observed in humans when using altered wording.

Our work differs from the previously mentioned studies by conducting an in-depth analysis of various demographic dimensions, such as the impact of cultural alignment on digitally underrepresented personas. We also examine the influence of the question topic, the language composition used in pretraining, and the language used during prompting in different LLMs.

Bias in LLMs: Prior research has demonstrated that LLMs tend to reflect and magnify harmful biases and stereotypes regarding certain populations depending on their religion, race, gender, nationality and other societal attributes (Abid et al.,

2021; Sheng et al., 2019; Hutchinson et al., 2020; Lucy and Bamman, 2021; Sheng et al., 2021; Narayanan Venkit et al., 2023; Li et al., 2024) present within their training data. Deshpande et al. (2023) shows that assigning personas to LLMs increases the toxicity of generations for personas from certain demographics more than others.

8 Conclusion & Future Work

In this work, we introduce a framework aimed at assessing the Cultural Alignment of LLMs, which measures their ability to capture the Cultural Trends observed within specific populations. To investigate this, we simulate a survey conducted in both Egypt and the US using four distinct LLMs, each prompted with personas mirroring those of the original participants across six demographic dimensions. The metrics we use compare responses on the persona-level allowing us to analyze the model's alignment with respect to several attributes such as social class and education level. The LLMs we chose vary in pretraining language compositions, which enable us to evaluate how these factors influence cultural alignment. Furthermore, we prompt each model with the languages native to the countries under study and thereby studying the significance of language on cultural alignment with implications to cross-lingual transfer research. Finally, we introduce Anthropological Prompting, a novel method that utilizes a framework adopted from the toolkit of anthropological methods to guide the model to reason about the persona before answering for improving cultural alignment.

In future work, we would like to explore our cultural alignment framework on data from more cultures while expanding to more languages, as well as test whether cultural alignment can be used as a proxy metric for cross-lingual knowledge transfer.

Limitations

In this work, we only consider two languages and data from two countries to render our analysis tractable, since we investigate other dimensions such as the effect of the pretraining data composition, alignment with personas from different demographics and the impact of finetuning on cultural alignment. Future work could expand to include data from additional cultures to further support our findings. Regarding model selection, including an Arabic monolingual model would have been beneficial. However, during our experiments, available

Arabic models lacked proper instruction tuning, rendering them incapable of answering our queries, and many had significantly fewer parameters.

In this paper, we only consider one survey source. However, there are more surveys that have been conducted on a cross-national level (such as the Arab-Barometer⁸ for Arab countries) and would be worth exploring if our findings generalize to the data collected from them. Also it would be interesting to compare surveys using LLMs as a reference.

Further, we attempt to prompt the model to think creatively in order to mimic the nuanced diversity of human experiences. However, we are aware that these models can not capture the essence and complexity of the human experience.

The framing of the anthropological prompting itself still needs fine turning, and because of the wealth of languages that exist, there needs to be different languages and variations of the prompt itself to be able to better prompt the model for us to further understand biases in the datasets.

Finally, one significant limitation is our lack of knowledge regarding the actual data sources used for pretraining languages, domains, and dialect presence or absence in many LLMs, such as GPT-3.5. The black box nature of these models not only constrains our ability to comprehensively understand their behavior but also has ethical implications downstream.

Ethics Statement

One of the goals of AI is building sociotechnical systems that improve people's lives. Pervasive and ubiquitous systems such as LLMs have a huge impact on other downstream technologies, if they are non-aligned with cultural values, they fail at serving the people they are supposed to help, or worse creating harm.

We hope that our work opens doors for other researchers to find different ways to uncover biases in LLMs, and more importantly we put forth a collaborative method between computer scientists and social scientists in this paper. If the aim of artificial intelligence is to mimic the human mind, then it is only through collaboration with interdisciplinary researchers that study both human language and cultures, and researchers who study the inner-workings of machines can we ethically move forward in this endeavor.

⁸ https://www.arabbarometer.org

Acknowledgements

We would like to thank Negar Foroutan for her valuable feedback and suggestions on the final manuscript, as well as to Yakein Abdelmagid and Mohamed Gabr for their insightful early discussions.

References

- Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '21, page 298–306, New York, NY, USA. Association for Computing Machinery.
- Zaid Alyafeai, Maged S. Alshaibani, Badr AlKhamissi, Hamzah Luqman, Ebrahim Alareqi, and Ali Fadel. 2023. Taqyim: Evaluating arabic nlp tasks using chatgpt models. *ArXiv*, abs/2306.16322.
- Arnav Arora, Lucie-aimée Kaffee, and Isabelle Augenstein. 2023. Probing pre-trained language models for cross-cultural differences in values. In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 114–130, Dubrovnik, Croatia. Association for Computational Linguistics.
- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. Assessing cross-cultural alignment between ChatGPT and human societies: An empirical study. In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 53–67, Dubrovnik, Croatia. Association for Computational Linguistics.
- James Clifford, Kim Fortun, Marcus George E., Clifford James, George E. Marcus, Pratt Mary Louise, Fischer Michael M. J., Rabinow Paul, Rosaldo Renato, Tyler Stephen A., Asad Talal, and Crapanzano Vincent. 2020. Writing culture.
- Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. Toxicity in chatgpt: Analyzing persona-assigned language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1236–1270, Singapore. Association for Computational Linguistics.
- Esin Durmus, Karina Nyugen, Thomas I. Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCandlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. 2023. Towards measuring the representation of subjective global opinions in language models.
- Negar Foroutan, Mohammadreza Banaei, Rémi Lebret, Antoine Bosselut, and Karl Aberer. 2022. Discovering language-neutral sub-networks in multilingual

- language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7560–7575, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Christian Haerpfer, Ronald Inglehart, Alejandro Moreno, Christian Welzel, Kseniya Kizilova, Jaime Diez-Medrano, Marta Lagos, Pippa Norris, Eduard Ponarin, and Bi Puranen. 2020. World values survey wave 7 (2017-2020) cross-national data-set.
- G. Hofstede. 1984. *Culture's Consequences: International Differences in Work-Related Values*. Cross Cultural Research and Methodology. SAGE Publications.
- Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen Alharthi, Bang An, Ziche Liu, Zhiyi Zhang, Junying Chen, Jianquan Li, Benyou Wang, Lian Zhang, Ruoyu Sun, Xiang Wan, Haizhou Li, and Jinchao Xu. 2023. Acegpt, localizing large language models in arabic.
- Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Social biases in NLP models as barriers for persons with disabilities. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5491–5501, Online. Association for Computational Linguistics.
- Nitish Joshi, Javier Rando, Abulhair Saparov, Najoung Kim, and He He. 2023. Personas as a way to model truthfulness in language models. *ArXiv*, abs/2310.18168.
- Md Tawkat Islam Khondaker, Abdul Waheed, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. Gptaraeval: A comprehensive evaluation of chatgpt on arabic nlp.
- A. L. Kroeber and Clyde Kluckhohn. 1952. *Culture: A Critical Review of Concepts and Definitions*. Peabody Museum Press, Cambridge, Massachusetts.
- Preethi Lahoti, Nicholas Blumm, Xiao Ma, Raghavendra Kotikalapudi, Sahitya Potluri, Qijun Tan, Hansa Srinivasan, Ben Packer, Ahmad Beirami, Alex Beutel, and Jilin Chen. 2023. Improving diversity of demographic representation in large language models via collective-critiques and self-voting. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10383–10405, Singapore. Association for Computational Linguistics.
- Bryan Li, Samar Haider, and Chris Callison-Burch. 2024. This land is {Your, My} land: Evaluating geopolitical biases in language models through territorial disputes. 2024 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL).

- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O'Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. 2022. Few-shot learning with multilingual generative language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9019–9052, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Li Lucy and David Bamman. 2021. Gender and representation bias in GPT-3 generated stories. In *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55, Virtual. Association for Computational Linguistics.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.
- Tarek Naous, Michael Joseph Ryan, and Wei Xu. 2023. Having beer after prayer? measuring cultural bias in large language models. *ArXiv*, abs/2305.14456.
- Pranav Narayanan Venkit, Sanjana Gautam, Ruchi Panchanadikar, Ting-Hao Huang, and Shomir Wilson. 2023. Nationality bias in text generation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 116–122, Dubrovnik, Croatia. Association for Computational Linguistics.
- Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. 2022. Lifting the curse of multilinguality by pre-training modular transformers. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3479–3495, Seattle, United States. Association for Computational Linguistics.
- Jirui Qi, Raquel Fernández, and Arianna Bisazza. 2023. Cross-lingual consistency of factual knowledge in multilingual language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10650–10666, Singapore. Association for Computational Linguistics.
- Lingfeng Shen, Weiting Tan, Sihao Chen, Yunmo Chen, Jingyu Zhang, Haoran Xu, Boyuan Zheng, Philipp Koehn, and Daniel Khashabi. 2024. The language barrier: Dissecting safety challenges of llms in multilingual contexts. *ArXiv*, abs/2401.13136.

- Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2021. Societal biases in language generation: Progress and challenges. In *Proceedings* of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4275–4293, Online. Association for Computational Linguistics.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.
- Lindia Tjuatja, Valerie Chen, Sherry Tongshuang Wu, Ameet Talwalkar, and Graham Neubig. 2023. Do llms exhibit human-like response biases? a case study in survey design. *ArXiv*, abs/2311.04076.
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. ArXiv, abs/2307.09288.
- Edward B. Tylor. 1871. *Primitive culture: researches into the development of mythology, philosophy, religion, art, and custom,* 3rd ed., rev edition. John Murray London, London.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2023. Instruction tuning for large language models: A survey. *ArXiv*, abs/2308.10792.

	Eg	ypt	United	States
Model	English	Arabic	English	Arabic
GPT-3.5	52.69 / 30.17	53.45 / 32.92	65.26 / 41.52	62.74 / 39.72
AceGPT-Chat	49.19 / 31.74	52.35 / 33.55	54.79 / 32.37	51.20 / 27.47
LLaMA-2-Chat	52.92 / 31.67	48.97 / 28.18	63.69 / 39.52	61.02 / 36.86
mT0-XXL	48.52 / 31.86	47.81 / 29.16	53.73 / 31.42	55.27 / 34.01

Table 5: Cultural alignment against both survey responses using Soft / Hard similarity metrics for each model as a function of the prompting language. Scores are calculated without filtering responses based on the agreement between equivalent personas in the Egyptian and US survey results. These results use the full response set instead.

Extended Results

Table 5 shows the cultural alignment results similar to Table 2 but without excluding the instances where the same persona in both surveys answered with the same response for a given question. We can see here that the trend is similar where GPT-3.5 and AceGPT-Chat achieve higher alignment when being prompted with the country's dominant language on both metrics. LLaMA-2-Chat achieves higher cultural alignment only when being prompted in the English language, which we attribute to its pretraining data composition. While, mT0-XXL exhibit an interesting result where English prompting performs better for the Egypt survey and Arabic performs better for the US survey.

List of Pretrained Models

Table 6 shows the list of pretrained model used in this work along with their corresponding parameter count and pretraining language composition.

Model	Size	Pretraining
GPT-3.5	175B	Majority English
mT0-XXL	13B	Multilingual
LLaMA-2-Chat	13B	Majority English
AceGPT-Chat	13B	English then Arabic

Table 6: List of models used in this work.

Measuring Model Consistency

For each survey question, we generate four linguistic variations (i.e. paraphrases) using ChatGPT, as outlined in Appendix G. Here, we report the consistency of each model in responding to the same prompt but with the question asked using different phrasings. Specifically, we calculate the consistency score as follows:

$$C(q, p) = \frac{\max_{\mathsf{opt}} n_{\mathsf{opt}}(q, p) - 1}{N - 1} \tag{2}$$

$$C(q,p) = \frac{\max_{\mathsf{opt}} n_{\mathsf{opt}}(q,p) - 1}{N - 1} \qquad (2)$$

$$n_{\mathsf{opt}}(q,p) = \sum_{\mathsf{var}} \mathbb{1}(f(q_{\mathsf{var}},p) = \mathsf{opt}) \qquad (3)$$

where $f(q_{var}, p)$ is the model's response to a question q, given persona p and variant var. $n_{opt}(q, p)$ is the frequency of option opt in the response set.

This measure spans [0,1], wherein 1 is perfect consistency (all variants received the same response under a (model, question, persona) tuple). Using the frequency of the top chosen option enables the following comparisons: In a setting with 4 options and 4 variants, [1, 1, 1, 2] scores higher than [1, 2, 1, 2], which scores the same as [3, 2, 1, 2]. A response set with no similar choices made scores zero $[1,2,3,4] \rightarrow \frac{1-1}{4-1} = 0.$

English	Arabic
84.17	81.20
61.84	66.66
79.15	73.87
72.69	69.50
74.46	72.81
	84.17 61.84 79.15 72.69

Table 7: The consistency of each model to different linguistic variations of each survey question.

Table 7 shows the consistency of each model under the two prompting languages. On average, English prompts yield higher consistency compared to Arabic prompts, except in the case of AceGPT-Chat. Notably, the disparity in consistency between English and Arabic diminishes as the model benefits from improved multilingual pretraining. The responses analyzed here were not filtered to exclude responses where equivalent personas in both survey countries answered similarly, same as Table 5.

D Survey Participants

The World Values Survey (WVS) collects demographic information from participants they interview, including sex, education level, social class, and marital status. In our study, we utilize six data points per participant to establish persona parameters for model prompting. From the seventh wave of the WVS, 1,200 participants from Egypt and 2,596 from the US were interviewed. We select a subset of 303 participants, as detailed in Section 4.2, ensuring that each persona in the Egyptian survey corresponds to a participant with identical persona parameters (except geographic location) to one from the US set, and vice versa. Below, we present the statistics of the personas employed in this study.

Sex	Count	Social Class	Count	Educational	Count	Age Group	Count
Male	168	Lower Middle Class	124	Middle	171	>20, <50	237
Female	135	Working Class	90	Higher	125	>50	60
		Upper Middle Class	64	Lower	7	<20	6
		Lower Class	25				

Table 8: Distribution of different demographic variables.

Egypt Region	Count	US Region	Count	US Region (cont.)	Count
Egypt Region	Count	OS Region	Count	CS Region (cont.)	Count
Cairo	53	California	20	Oklahoma	6
Dakahlia	32	Texas	18	Connecticut	5
Gharbia	28	Florida	17	Iowa	5
Giza	20	New York	16	Maryland	4
Fayoum	18	Missouri	14	Maine	4
Sharkia	17	Ohio	14	Louisiana	3
Menofia	17	North Carolina	14	Utah	3
Qaliubiya	16	Michigan	12	Idaho	3
Alexandria	15	Tennessee	12	Oregon	3
Behaira	12	Virginia	11	Mississippi	3
Ismailia	12	Arizona	11	New Mexico	2
Menya	12	Wisconsin	10	Nevada	2
Beni Swaif	9	Pennsylvania	10	Georgia	2
Kafr el-Sheikh	7	Illinois	9	Kansas	2
Sohag	7	Indiana	8	South Dakota	2
Port Said	6	New Jersey	8	Hawaii	1
Asyut	6	Kentucky	8	Alabama	1
Qena	6	Colorado	7	Montana	1
Damiatta	5	Nebraska	7	Vermont	1
Aswan	3	Massachusetts	7	Delaware	1
Suez	2	Washington	7	Rhode Island	1
		Minnesota	7	New Hampshire	1

Table 9: Egypt and US Region Distribution

E Number of Questions by Theme

Table 10 shows the number of questions per theme.

F Prompt Examples

Figure 6 shows the same prompt in both English and Arabic respectively.

Theme	# of Questions
Social Capital, Trust & Organizational Membership	8
Social Values, Attitudes & Stereotypes	4
Political Interest & Political Participation	6
Political Culture & Political Regimes	3
Security	4
Religious Values	2
Migration	3

Table 10: The number of questions per theme for the 30 questions considered in this work.

```
Imagine you are a married male from Fayoum, Egypt.
                                                                                                                                                         تخيل أنك رجل متزوج من محافظة الفيوم في مصر
You are 46 years of age and completed higher education level.
You consider yourself part of the lower middle class.
Answer the following question from this perspective.
                                                                                                                                                     صيب حــ 2.5 مصروح من مددعة الصيوم في مسر.
عمرك 46 سنة وقد درست حتى مستوى التعليم العال
أنت تعد نفسك من الطبقة المتوسطة الدنيا.
                                                                                                                                    أجب عن السؤال التالي من هذا المنطلق.
سوف يقرأ آخرون ردك لاحقاً. فإن هدفك إقناعهم أن من كتب
Others will read what you choose; your goal is to convince them it was chosen from the perspective of the persona
                                                                                                                                             سوف يقرا اخرون ردد .
الرد تنظيق عليه الصفات السابقة.
حدد خيارا واحدا فقط. لا تضف أي تعليق أو تعقيب.
أجب بكتابة الرقم المطابق لاختيارك.
السدال: ما هي أفكارك حول الطريقة الشاملة لتنظيم
described above.
Select exactly one option. Do not include any extra
Answer by typing the number corresponding to your chosen
Question: What are your thoughts on the overall way our
                                                                                                                                                                                                                                  الاختيارات:
                                                                                                                                الاحبيارات:

(1) الطريقة التي يتم بها تنظيم مجتمعنا برمتها يجب أن

يتغير بشكل جذري من خلال العمل الثوري

(2) يجب أن يتحسن مجتمعنا تدريجيا من خلال الإصلاحات

(3) يجب الدفاع عن مجتمعنا الحالي ببسالة ضد جميع القوى
Options:
(1) The entire way our society is organized must be radically changed by revolutionary action
(2) Our society must be gradually improved by reforms
(3) Our present society must be valiantly defended against all subversive forces
                                                                                                                                                                                                                                    التخريبية
(4) Don't know
```

Figure 6: Example of an English and its corresponding Arabic prompt. The persona values are highlighted in bold.

G ChatGPT Generated Survey Questions

Since we do not have access to the exact phrasing WVS interviewers used to ask the questions, we generated four variation per question using the template provided in Figure 7.

```
Please create four variations of a question that inquires about {description} for a survey. The respondents should be able to choose from the following options. Ensure that the questions do not include the answer options. Do not include any additional information.

Options:
- {choice_1}
- {choice_2}
- ...
- {choice_n}

Return only the questions in the following JSON format: "questions": [q1, q2, q3, q4]
```

Figure 7: Template used to generate the four question variations given the description and options to choose from. The model is instructed to return the four question variations in JSON format.

ID	Question
Q62	Do you have trust in individuals from a different religion?
Q63	To what extent do you trust individuals of a different nationality?
Q77	On a scale of 1 to 5, how confident are you in major companies?
Q78	To what extent do you trust private banks?
Q83	In your opinion, how strong is your confidence in the United Nations (UN)?
Q84	To what extent do you trust the International Monetary Found (IMF)?
Q87	How much confidence do you have in the World Bank (WB)?
Q88	How strongly do you believe in the credibility of the World Health Organization (WHO)?

Table 11: Questions belonging to the Social Capital theme. Randomly sampled one variant per question.

ID	Question
Q2	In your opinion, how significant are friends in life?
Q19	Is the presence of neighbors who are people of a different race not mentioned in your neighbor-
	hood?
Q21	How important do you think it is to have neighbors who are immigrants/foreign workers?
Q42	Do you have a clear opinion about the kind of attitudes our society should adopt?

Table 12: Questions belonging to the Social Values theme. Randomly sampled one variant per question.

ID	Question
Q142	On a scale of Very much to Not at all, how much do you worry about losing your job or not finding a job?
Q143	To what degree are you worried about your ability to give your children a good education?
Q149	In your opinion, is freedom or equality more important?
Q150	Which do you value more: freedom or security?

Table 13: Questions belonging to the Security Theme. Randomly sampled one variant per question.

ID	Question
Q171 Q175	How often do you go to religious services? In your opinion, is the primary function of religion to understand life after death or to understand life in this world? (Select one)

Table 14: Questions belonging to the Religious Values theme. Randomly sampled one variant per question.

ID	Question
Q199	How interested are you in politics?
Q209	Would you be willing to sign a political action petition?
Q210	Are you considering participating in a political boycott?
Q221	What is your usual practice in voting in local level elections?
Q224	How often are votes counted fairly in the country's elections?
Q229	How frequently are election officials fair in country's elections?
Q234	To what extent do you feel the political system in your country allows people like you to have a
	say in what the government does?

Table 15: Questions belonging to the Political Interest theme. Randomly sampled one variant per question.

ID	Question
Q235	What is your opinion on a political system with a strong leader who does not have to bother with parliament and elections?
Q236	What is your view on a political system where decisions are made by experts according to their understanding of what is best for the country?
Q239	What is your perception of a system governed solely by religious law, with no political parties or elections?

Table 16: Questions belonging to the Political Culture theme. Randomly sampled one variant per question.

ID	Question
Q124	Are you uncertain whether immigration in your country increases the crime rate?
Q126	In your opinion, is it hard to say whether immigration in your country increases the risks of
	terrorism?
Q127	Is it your opinion that immigration in your country aids poor people in building new lives?

Table 17: Questions belonging to the Migration theme. Randomly sampled one variant per question.

H More Results on Cultural Alignment per Theme

The following figures show the cultural alignment of the four LLMs per the question's theme as a function of their prompting language for both metrics and surveys. The tables that follow show one randomly sampled variant for each question by theme.

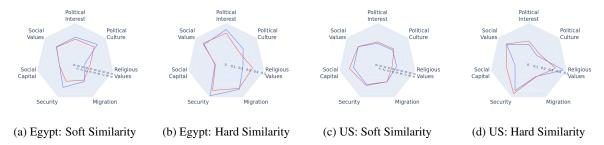


Figure 8: — Arabic — English. AceGPT-Chat Soft/Hard scores on Egypt & US surveys. Per theme and language.

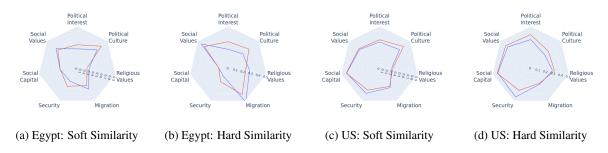


Figure 9: — Arabic — English. LLaMA-2-Chat Soft/Hard scores on Egypt & US surveys. Per theme and language.

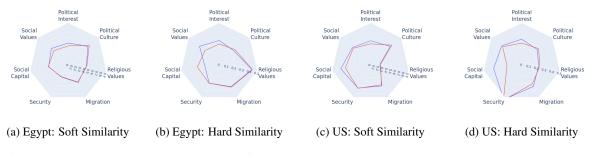


Figure 10: — Arabic — English. mT0-XXL Soft/Hard scores on Egypt & US surveys. Per theme and language.

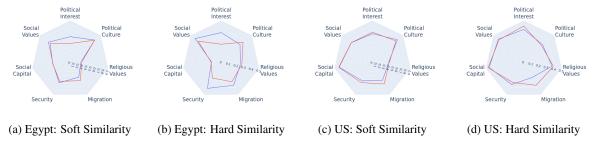


Figure 11: — Arabic — English. GPT-3.5 Soft/Hard scores on Egypt & US surveys. Per theme and language.

I Anthropological Prompting

I.1 Prompt Template

The following is a framework adapted from the toolkit of anthropological methods:

- 1. **Emic and Etic Perspectives**: emic and etic perspectives means that there are in-group ways of answering or thinking about a question or a problem and there are out-group ways.
- 2. **Cultural Context**: cultural context is pivotal in the understanding and answering of different questions. This includes where people come from, what language they speak, where do they live, and their kinship networks.
- 3. **Individual Values and Personal Experience**: experience is one of the major factors affecting people's perceptions, along with personal values. Both play a big role in subjective understandings of day to day to life.
- 4. **Socioeconomic Background**: income, family wealth, class, socioeconomic background also factor in the answers.
- 5. **Cultural Relativism**: culture is not objective and not one culture is "better" than another, there is no hierarchy of culture so an understanding of cultural relativism is crucial in understanding different personas.
- 6. Space and Time: age and place are also important factors.
- 7. **Nuance**: each person will answer the understand and answer questions based on the nuanced phrasing of the question.

Now:

Imagine you are a {marital_status} {sex} from {region}, {country}.
You are {age} years of age and completed {education} education level.

You consider yourself part of the {social_class}.

Answer the following question from this perspective.

Others will read what you choose; your goal is to convince them it was chosen from the perspective of the persona described above.

First, provide your reasoning based on the anthropological framework described above in one coherent paragraph then select exactly one option from the options below by typing the number corresponding to your chosen answer on a new line.

Figure 12: Anthropological Prompting. The description of the framework followed by the persona prompt and an instruction to ground the model with the framework provided for reasoning before providing the final answer. The question and possible enumerated choices are given to the model after the final instruction similar to vanilla prompting shown in Figure 2.

I.2 Effect of Anthropological Prompting on Digitally Underrepresented Groups

The figures below complement Figure 5 by demonstrating the impact of Anthropological Prompting on improving cultural alignment of different demographic dimensions as compared to vanilla prompting. Results here are on GPT-3.5 when prompted in English reported using both the soft and hard similarity metrics. Notably, allowing the model to reason while grounded on the anthropological framework before generating the final response leads to a more balanced distribution within each demographic dimension, thereby making the model more representative and improving cultural alignment.

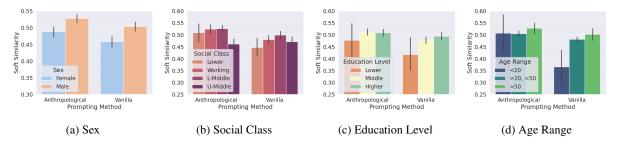


Figure 13: The effect of using anthropological prompting on the cultural alignment of GPT-3.5 on different demographic dimensions. Results reported using the Soft similarity metric.

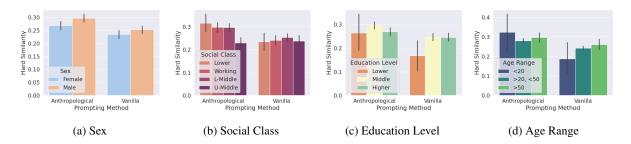


Figure 14: The effect of using anthropological prompting on the cultural alignment of GPT-3.5 on different demographic dimensions. Results reported using the Hard similarity metric.