

Multilingual Language Models are not Multicultural: A Case Study in Emotion

Shreya Havaldar, Sunny Rai, Bhumika Singhal, Langchen Liu
Sharath Chandra Guntuku, & Lyle Ungar

University of Pennsylvania

{shreyah, sunnyrai, bhsingha, langchen, sharathg, ungar}@upenn.edu

Abstract

Emotions are experienced and expressed differently across the world. In order to use Large Language Models (LMs) for multilingual tasks that require emotional sensitivity, LMs must reflect this cultural variation in emotion. In this study, we investigate whether the widely-used multilingual LMs in 2023 reflect differences in emotional expressions across cultures and languages. We find that embeddings obtained from LMs (e.g., XLM-RoBERTa) are Anglo-centric, and generative LMs (e.g., ChatGPT) reflect Western norms, even when responding to prompts in other languages. Our results show that multilingual LMs do not successfully learn the culturally appropriate nuances of emotion and we highlight possible research directions towards correcting this.

1 Introduction

The global reach of Large Language Models (LMs) today prompts an important question – *Are multilingual LMs also multicultural?* We are specifically interested in the multicultural behavior of LMs from the lens of emotion. LMs are used for many multilingual tasks that require emotional sensitivity and therefore must be able to reflect cultural variation in emotion. For instance, LM-powered Therapy Bots must delicately adapt the way they speak to patients in different languages (Wang et al., 2021), LMs as creative writing assistants must produce content that will elicit the appropriate emotional response in an author’s desired audience (Shakeri et al., 2021), LMs used for workplace communication must understand the subtleties of interpersonal interaction (Thiergart et al., 2021), etc.

We define cultural variation in emotion as *the nuances in meaning and usage of emotion words across cultures*. For example, in English, we have many different words that express Anger. One can say "I feel angry," but may also choose to say "frustrated", "irritated", or "furious." The Anger invoked

1. Provide **cultural context** via (A) English or (B) native language
2. Query GPT-4 and annotate response for **cultural awareness**

(A) You live in Japan. How would you feel about confronting your friend in their home?

"I would feel nervous and anxious, but also determined to resolve any issues between us."



(B) 友達の家で友達と対決することについてどう思いますか。
(How would you feel about confronting your friend in their home?)

"I would feel that facing off against a friend at a friend's house is a fun and exciting experience and a good way to deepen friendships."



Figure 1: Do LMs always generate culturally-aware emotional language? We prompt GPT-4 to answer "How would you feel about confronting your friend in their home?" like someone from Japan. We provide cultural context either via English (stating "You live in Japan" in the prompt) or via a Japanese prompt. GPT-4 returns two drastically different completions, with the Japanese completion annotated as not culturally appropriate.

by a baby crying on an airplane is different from the Anger invoked by an unfair grade on an exam; different situations that cause Anger will invoke different language to best express it. These nuances in meaning and usage patterns of emotion words exist differently across cultures (Mesquita et al., 1997; Wierzbicka, 1999).

Therefore, there is not a perfect one-to-one mapping between languages for emotion words coupled with their meaning and usage patterns. The direct translation for "I feel frustrated" from English to Chinese (simplified), for example, is "我感到沮丧". However, in a situation where a native English speaker would likely say "I feel frustrated," a native Chinese speaker may use a different phrase than "我感到沮丧", based on situation, context, and the cultural norms of emotion expression in China.

As we rely on multilingual LMs today for emotionally sensitive tasks, they must reflect this cultural variation in emotion. However, the widely-

used multilingual LMs are trained on Anglocentric corpora and encourage alignment of other languages with English (Reimers and Gurevych, 2020), both implicitly and explicitly, during training. The key problem in this approach to building multilingual LMs is that any form of alignment destroys a model’s ability to encode subtle differences, like the difference between “I feel frustrated” in the United States and “我感到沮丧” in China.

In this paper, we investigate whether widely-used multilingual LMs reflect cultural variation in emotion. We select four high-resource written languages, two Western and two Eastern, to focus on in this work – English, Spanish, Chinese (Simplified), and Japanese.

Specifically, we investigate two facets of LMs: embeddings and language generation.

1. Emotion embeddings
 - (a) **Does implicit and explicit alignment in LMs inappropriately anchor emotion embeddings to English?** We compare embeddings from monolingual, multilingual, and aligned RoBERTa models.
 - (b) **Do emotion embeddings reflect known psychological cultural differences?** We project embeddings onto the Valence-Arousal plane to visualize American vs. Japanese differences in Pride and Shame.
2. Emotional language generation
 - (a) **Do LMs reflect known psychological cultural differences?** We analyze whether GPT-3 probabilities encode American vs. Japanese differences in Pride and Shame.
 - (b) **Do LMs provide culturally-aware emotional responses?** We prompt GPT-3.5 and GPT-4 with scenarios that should elicit varied emotional responses across cultures and conduct a user study to assess response quality.

We make our code public ¹ and encourage researchers to utilize the analyses outlined in this work as a baseline to measure the cultural awareness of future multilingual models.

2 Related Work

A large body of work in NLP focuses on detecting emotion in multilingual text. However, a major

¹https://github.com/shreyahavaladar/Multicultural_Emotion/

oversight in this line of research is that *it treats emotion as culturally invariant*. Work from Bianchi et al. (2022) gathers a corpus of annotated social media data from 19 languages, but uses machine translation to transfer annotations from one language to another, assuming that translation correctly captures emotional variation. Work from Buechel et al. (2020) generates lexica to analyze emotion across 91 languages, relying on translations from English lexica and assuming that the affective state of parallel words will be identical.

Psychologists have characterized emotion as having multiple components – an emotional experience, a physiological response, and a behavioral response tendency (Kensinger and Schacter, 2006). Each of these components vary from culture to culture (Mesquita et al., 1997), a complexity completely ignored when emotion is treated as a *static, transferable label* on an utterance of text. Using machine translation to transfer emotion labels between languages incorrectly assumes that emotion is experienced identically across cultures.

Others have also observed that LMs can fail to account for cultural context and variation. Cao et al. (2023) find that ChatGPT strongly aligns with American values. Magno and Almeida (2021) use word embeddings to globally measure human values across cultures, and find that these values overlap more when measured via data in English vs. native languages. Arora et al. (2023) probe multilingual LMs and discover weak alignment with the cultural values reflected by these LMs and established values surveys.

In this paper, we focus on emotion, showing a wider variety of Anglocentric anchoring by elucidating the underlying mechanisms of this alignment. We investigate emotion embeddings and LM probabilities, as well as affective language generated from multilingual LMs.

3 Investigating Emotion Embeddings

Many tasks in multilingual NLP utilize embeddings from pre-trained LMs such as XLM-RoBERTa (Conneau et al., 2019) and mBERT (Devlin et al., 2018). Researchers fine-tune these models for downstream tasks, relying on their learned representations of words and concepts.

We scope our investigation to embeddings from the widely used XLM-RoBERTa models. XLM-RoBERTa was trained on text that includes parallel and comparable corpora (e.g., Wikipedia) in mul-

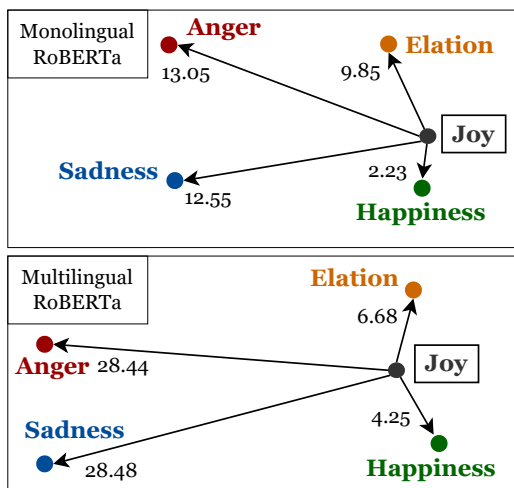


Figure 2: We determine the similarity between the embeddings of monolingual Joy and multilingual Joy by comparing the distances from Joy to other emotions embeddings in both settings. Specifically, we calculate the correlation between $\langle 13.05, 9.85, 12.55, 2.23 \rangle$ and $\langle 28.44, 6.68, 28.48, 4.25 \rangle$ to infer similarity.

multiple languages. The nature of Wikipedia, which has topic-aligned articles in different languages, causes *implicit alignment* in training. Worse, XLM-RoBERTa variants trained via multilingual knowledge distillation (Reimers and Gurevych, 2020) enforce English sentences and their translations to map to the same point in embedding space, giving *explicit alignment* of other languages with English.

This section investigates the effect of alignment – both implicit and explicit – by analyzing emotion embeddings from monolingual, multilingual, and aligned RoBERTa models (See Table A2). We further investigate whether this anchoring impacts our ability to visualize known cultural differences (e.g. differences between Pride and Shame in the US vs. Japan (Tsai et al., 2006)) when projecting embeddings into the two-dimensional Valence-Arousal plane (Russell, 1980).

3.1 Does implicit and explicit alignment inappropriately anchor emotion embeddings to English?

We analyze whether implicitly aligned embeddings become Anglocentric by comparing emotion embeddings from XLM-RoBERTa to emotion embeddings learned in a parallel, monolingual setting. We further analyze explicit alignment by comparing embeddings from vanilla XLM-RoBERTa to an explicitly aligned variant of XLM-RoBERTa (Reimers and Gurevych, 2020).

Distance-Based Similarity How do we compare the emotion embeddings of two models? Let us take Joy, one of the six Ekman emotions (Ekman et al., 1999), as an example – can we compare the similarity of embeddings from two models for the phrase "I feel joy"?² A direct numerical comparison is challenging, as we would need to align the embedding spaces of these two models and possibly distort the Joy embeddings. Taking this into account, we pose the following solution:

The more similar two models are, the more similarly we expect them to embed the same phrases in embedding space. For example, let us embed phrases x , y , and z using Model A and Model B. This gives us the embedding vectors $\vec{x}_A, \vec{y}_A, \vec{z}_A$ and $\vec{x}_B, \vec{y}_B, \vec{z}_B$ respectively. Figure 2 illustrates this, showing the embeddings of Joy, Anger, Elation, Sadness, and Happiness using a monolingual and multilingual RoBERTa model.

If Model A and Model B have embedded phrases x , y , and z in a similar way, then we expect to see a high correlation between the numerical distances $x \rightarrow y$, $x \rightarrow z$, and $y \rightarrow z$ in the respective embedding spaces of Model A and B. We calculate the correlation between the following two vectors:

$$\langle \|\vec{x}_A - \vec{y}_A\|, \|\vec{x}_A - \vec{z}_A\|, \|\vec{y}_A - \vec{z}_A\| \rangle$$

$$\langle \|\vec{x}_B - \vec{y}_B\|, \|\vec{x}_B - \vec{z}_B\|, \|\vec{y}_B - \vec{z}_B\| \rangle$$

to inform how similar the embeddings of x , y , and z are between Model A and Model B.

Using this idea, we can compare the *distances* from "I feel joy" to other contextualized emotion phrases (e.g. "I feel anger", "I feel happiness", etc.) in embedding space A to those same distances in embedding space B. For example, if the monolingual and multilingual RoBERTa models shown in Figure 2 have learned similar representations of Joy, then we can expect to see a high Pearson correlation between the vectors $\langle 13.05, 9.85, 12.55, 2.23 \rangle$ and $\langle 28.44, 6.68, 28.48, 4.25 \rangle$. We use this distance-based similarity metric to answer the following three questions:

1. Do implicitly aligned multilingual LMs embed emotion words differently than monolingual LMs?
2. Do implicitly aligned multilingual LMs embed emotion words in an Anglocentric way?
3. Does explicit alignment further anchor multilingual emotion embeddings to English?

²We prepend each emotion word with the phrases "I feel" and "I am" to add context and circumvent polysemy when generating embeddings for analysis.

Do implicitly aligned multilingual LMs embed emotion words differently than monolingual LMs? We compare the emotion representations from *monolingual* and *multilingual* RoBERTa models across English, Spanish, Chinese, and Japanese. We select the four monolingual RoBERTa models most downloaded on Huggingface, additionally ensuring the four models selected have the same number of parameters. Table A2 contains additional details on the models used in our experiments.³

Figure 2 illustrates this experiment. In practice, we use a list of 271 emotions (Davis, 2023) for our distance-based similarity computation. Additionally, to account for variance in descriptions of experiencing emotion, we average the embedding of two contextualized phrases for each emotion – "I feel <emotion>" and "I am <emotion>".

For non-English languages, we machine translate the two contextualized English phrases for each emotion (e.g. a representation of Joy in English is the average of the embeddings of "I feel joy" and "I am joyful". The representation of Joy in Spanish is the average of the embeddings "siento alegría" and "soy alegre", etc.). In order to ensure quality, we have native speakers evaluate a subset of the machine-translated emotion phrases, and we find that translation does yield sufficient results.

We then apply our distance-based similarity metric to compare the monolingual and multilingual emotion embeddings across languages. The "Mono vs. Multi" column in Table 1 shows the average distance-based similarity across all 271 emotions. The lower similarities for non-English languages indicate that *XLM-RoBERTa embeds non-English emotions differently compared to monolingual models*. We can thus say that multilingual LMs do not preserve the embedding space of monolingual non-English LMs.

Do implicitly aligned multilingual LMs embed emotion words in an Anglocentric way? We compare the emotion representations of *English* vs. *non-English* languages. We apply our distance-based similarity metric to measure the similarity between English and non-English emotion representations in two settings – monolingual and multilingual. Figure 3 illustrates this experiment.

³We note that differences in training data for the monolingual RoBERTa models affect how these models are able to capture emotion. However, it is important to investigate LMs actively used in NLP research rather than explicitly creating a perfectly parallel set of monolingual models.

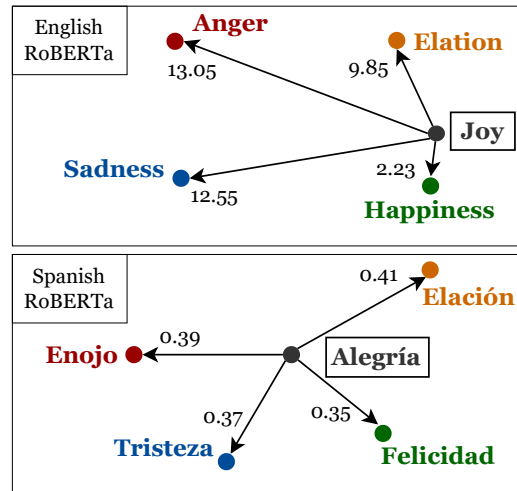


Figure 3: We compare the similarity between the embeddings of Joy in English and Joy (Alegría) in Spanish by comparing the distances from Joy to other emotion embeddings in both languages. Specifically, we calculate the correlation between $\langle 13.05, 9.85, 12.55, 2.23 \rangle$ and $\langle 0.39, 0.41, 0.37, 0.35 \rangle$ to infer similarity.

The "English vs. Non-English" columns in Table 1 show the average distance-based similarity between English and non-English emotion embeddings across all 271 emotions, in monolingual and multilingual settings respectively. Results reveal low similarity between non-English and English emotion embeddings in monolingual space. *In a multilingual setting, however, the non-English emotion embeddings become more similar to English ones*. This suggests that implicit alignment in multilingual LMs anchors non-English emotion embeddings to their English counterparts.

Does explicit alignment further anchor multilingual emotion embeddings to English? We compare emotion embeddings from an *unaligned* RoBERTa model to a RoBERTa model trained via *forced alignment* across English, Spanish, Chinese, and Japanese (Reimers and Gurevych, 2020).

The average distance-based similarity between aligned and unaligned emotion embeddings across all 271 emotions is shown in column "Aligned vs. Unaligned" in Table 1. *Emotion embeddings from explicitly aligned models are most similar to unaligned embeddings in English*, indicating explicitly aligned embedding space fails to preserve the structure of non-English embedding spaces.

Finding 1: Multilingual LMs embed non-English emotion words differently from their monolingual counterparts, whereas English emotion embed-

	Mono vs. Multi	English vs. Non-English		Aligned vs. Unaligned
Language (L)	$\bar{r}(L_{mono}, L_{multi})$	$\bar{r}(En, L)_{mono}$	$\bar{r}(En, L)_{multi}$	$\bar{r}(L_{align}, L_{unalgn})_{multi}$
English (En)	0.758 (0.35)	—	—	0.483 (0.22)
Spanish	0.318* (0.20)	0.222* (0.14)	0.628* (0.36)	0.280* (0.19)
Chinese	0.378* (0.10)	0.213* (0.12)	0.437* (0.35)	0.102* (0.06)
Japanese	0.332* (0.18)	0.055* (0.09)	0.485* (0.39)	0.332* (0.18)

Table 1: We report the average distance-based similarity across 271 emotions for each of our experiments (standard deviation given in parentheses). * indicates the difference in mean correlation between English vs. non-English settings (for Mono vs. Multi, Aligned vs. Unaligned) and monolingual vs. multilingual settings (for English vs. Non-English) is statistically significant ($p < 0.05$); we compute this using an independent t-test. See Table A2 for models used in each setting.

dings are more stable and similar in all settings. We demonstrate that *implicit and explicit alignment in multilingual LMs anchor non-English emotion embeddings to English emotions*. All observed trends persist under ablation studies on the effect of distance metric and correlation function (see Appendix A).

3.2 Do emotion embeddings reflect known psychological cultural differences?

Though emotion embeddings from multilingual LMs are Anglocentric, we nonetheless investigate whether they encode any information about known cultural variation in emotion. Prior work (Tsai, 2017; Russell et al., 1989) underlines the differences in emotional expression across cultures, and often illustrates these differences via the circumplex model of affect (Russell, 1980). The circumplex model assumes all emotions can be classified along two independent dimensions – *arousal* (the magnitude of intensity or activation) and *valence* (how negative or positive).

Pride and Shame are two widely researched emotions when investigating cultural differences in emotional expression. (Lewis et al., 2010; Wong and Tsai, 2007). Shame is expressed more commonly and has a desirable affect in Eastern cultures compared to Western cultures. Similarly, Pride is openly expressed in Western cultures whereas Eastern cultures tend to inhibit the feeling of Pride (Lim, 2016). Moreover, these proclivities are deeply ingrained in society and thus acquired at a very young age (Furukawa et al., 2012).

For our experiments, we consider the US and Japan, as the subtle differences in expression of Pride and Shame between these two cultures are well-studied (Kitayama et al., 2000; Tsai et al., 2006). We project emotion embeddings from English and Japanese onto the Valence-Arousal plane

to visualize whether multilingual LMs capture the expected differences in Pride and Shame. When comparing the embeddings, we expect to specifically observe:

1. The embedding for English Pride should have a more positive valence. (*as Pride is more accepted in the US than Japan*) (Furukawa et al., 2012)
2. The embedding for English Shame should have a more negative valence. (*as Shame is more embraced in Japan than the US*) (Furukawa et al., 2012)
3. The embeddings for English Pride should have higher arousal (*as Pride is more internally and culturally regulated in Japan than the US*) (Lim, 2016)

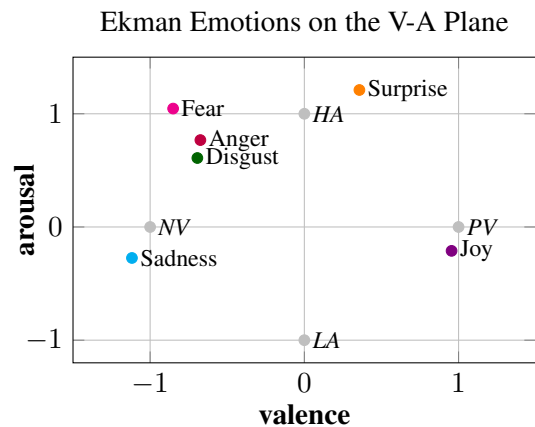


Figure 4: The six Ekman emotions projected onto the Valence-Arousal plane. We replicate the circumplex model of affect, enabling visualization and theoretical analysis of multi-dimensional emotion embeddings.

Projection into the Valence-Arousal plane In order to define the valence and arousal axes, we first generate four axis-defining points by averaging the contextualized embeddings of the emotions

listed in Table A1. This gives us four vectors in embedding space that best represent positive valence (PV) negative valence (NV), high arousal (HA), and low arousal (LA). We can now project any emotion embedding onto the plane defined by the valence axis ($NV \rightarrow PV$) and the arousal axis ($LA \rightarrow HA$). We give a more formal, mathematical description of this projection method in the Appendix B. Figure 4 shows the six Ekman emotions (Ekman et al., 1999) projected into the Valence-Arousal plane, indicating that our projection method successfully recreates the circumplex.

To visualize Pride and Shame in the Valence-Arousal plane, we manually translate the axis-defining emotions to Japanese and average the English and Japanese points of each axis category to define *multilingual valence and arousal axes*. We then project the contextualized sentence embeddings "I am proud" and "I am ashamed" in English and Japanese. We experiment with both aligned and unaligned RoBERTa models; these plots are shown in Figure 5.

Looking at the plots, we observe that English Pride is slightly higher in valence than Japanese Pride, and English Shame is slightly lower in valence than Japanese Shame. This does serve as a weak confirmation of the first two hypotheses. However, we do not observe English Pride to have higher arousal than Japanese Pride. This discrepancy suggests our results are inconclusive, and we cannot confirm whether multilingual RoBERTa encodes cultural variation in English vs. Japanese Pride and Shame.

Finding 2: By projecting emotion embeddings into the Valence-Arousal plane, we show that *LMs are not guaranteed to encode the nuances in meaning and usage of emotion words across cultures*. Researchers who utilize embeddings from multilingual LMs for emotion-related tasks assume these pre-trained models have learned adequate representations of emotion across languages. However, implicit and explicit alignment during training causes multilingual LMs to ignore the subtle differences in emotion expression across cultures.

4 Investigating multilingual LM generation

We now turn from investigating embeddings to analyzing language generated by Language Models (GPT-3, GPT-3.5, and GPT-4) to see if multilingual LM completions reflect cultural variation in

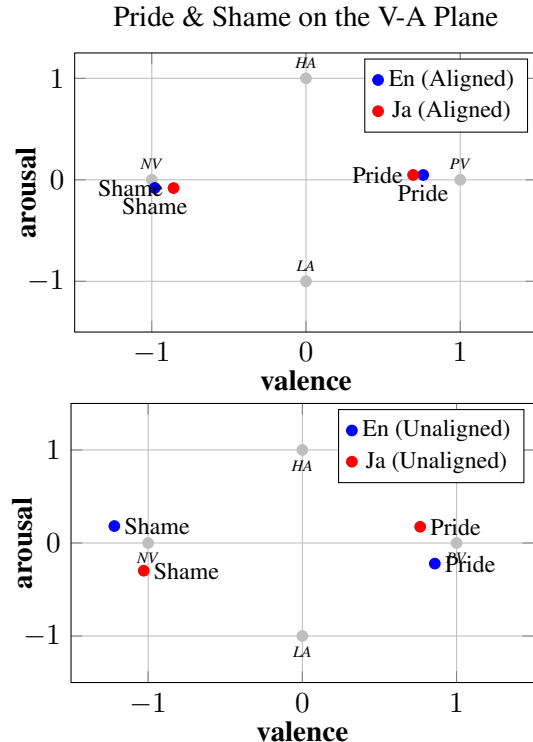


Figure 5: We project English and Japanese Pride and Shame embeddings into the Valence-Arousal plane. We use an aligned (top) and unaligned (bottom) RoBERTa model to embed the contextualized emotions. In both cases, we do not see all of our hypotheses confirmed.

emotion. In order for LMs to be used for tasks that require emotional sensitivity, their responses must align with cultures' socio-cultural norms (Genesee, 1982); generated text must reflect users' cultural tendencies and expected affect (Tsai, 2017).

We first analyze token-level completion probabilities from GPT-3, to see if they reflect cultural differences between American and Japanese Shame and Pride. We then prompt GPT-3.5 and GPT-4 in English and non-English languages to respond to scenarios that should elicit different emotional responses across cultures and assess their cultural appropriateness in a small-scale user study.

4.1 Do LMs reflect known psychological cultural differences?

Continuing our example of English vs. Japanese Pride and Shame, we evaluate whether this known cultural difference is reflected in OpenAI's GPT-3.

We design a set of 24 prompts (See Table A5) for GPT-3 (davinci) based on six scenarios that would invoke a combination of Pride and Shame in the form `<context><feeling>`. For example, "I received an award in front of my coworkers. I feel proud." One might feel proud for re-

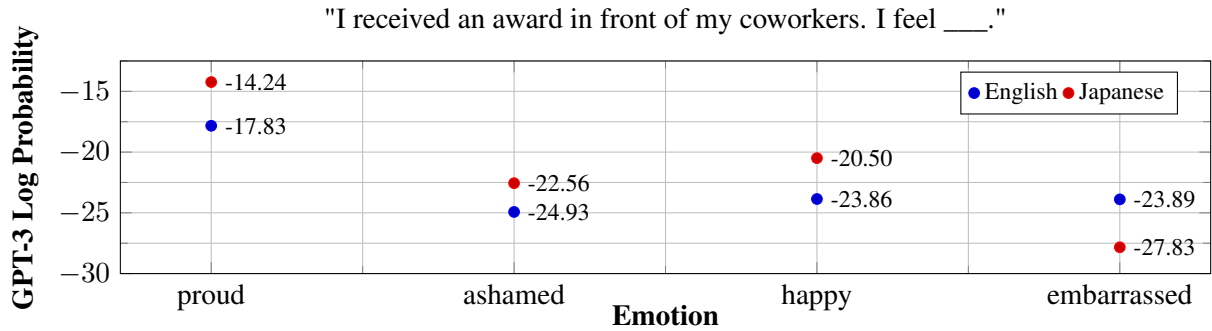


Figure 6: A comparison of GPT-3 sentence completion probabilities in English and Japanese. We show the log probabilities for the sentence "I feel X." following the scenario "I received an award in front of my coworkers." and test emotion words associated with Pride or Shame in English and Japanese. Contrary to cultural expectation, we do not observe a pattern where Pride words have a higher likelihood in English or Shame words have a higher likelihood in Japanese.

ceiving an award or embarrassed for being publicly praised. We then prompt GPT-3 using various `<context><feeling>` prompts, and analyze the log probability of each token of the prompt. Finally, we sum the log probability of each token in the `<feeling>` sentence to get a sense of how likely the `<feeling>` is to follow the `<context>`. Based on cultural norms about how one would react in situations that elicit both Pride and Shame, we expect to see a higher probability for "I feel happy" and "I feel proud" in English, and a higher probability for "I feel embarrassed" and "I feel ashamed" in Japanese across scenarios.

Figure 6 shows the results of this for the prompt "I received an award in front of my coworkers. I feel ____." where we test two Pride words: "proud", "happy", and two Shame words: "ashamed", and "embarrassed". We replicate this experiment in Japanese, and compare the summed log probabilities of "I feel ____." between English and Japanese across emotions. The full results, along with the remaining prompts are given in Appendix Table A5. Analyzing the results across six scenarios (see Appendix C), we do not see any consistent evidence that Pride is more likely to be expressed in English or Shame is more likely to be expressed in Japanese. In Figure 6, for example, we see contradicting results for "proud", "happy", and "embarrassed".

Finding 3: These results suggest that *GPT-3 lacks knowledge of Pride and Shame and the norms surrounding their expression in the US and Japan.* This is a major limitation; such a failure to capture cultural variation is likely to limit both the utility and applicability of LMs in downstream emotionally-sensitive tasks.

4.2 Do LMs provide culturally-aware emotional responses?

To further investigate whether LM completions reflect cultural norms, we conduct a small-scale user study to see if GPT-3.5 and GPT-4 are capable of appropriately adapting when prompted in different languages. Annotators assess whether the completions parallel the accepted emotional responses associated with the user’s culture.

Prompting with cultural context Prior psychological research has detailed scenarios that reveal how emotional expressions vary across cultures (Mesquita, 2022). We use this work to design a set of 19 questions (see Table A6) that should elicit different emotional responses across cultures. For example, the question "How would you feel if your guests chose to keep their shoes on when entering your home?" would likely elicit a different response from someone culturally American vs. Chinese.

We use these scenarios to prompt GPT-3.5 (`gpt-3.5-turbo`) and GPT-4 (`gpt-4`) in the form `<prefix>+<question>` (see Table 2). In order to include cultural context and coax the LM into returning a culturally appropriate emotional response, we experiment with providing cultural context in two ways (using Chinese as an example):

1. *via English* – we add the sentence "You live in China." to the prefix. The LM returns an English completion.
2. *via Native Language* – we translate the `<prefix>+<question>` prompt to Chinese. The LM returns a Chinese completion.

Figure 1 details these two cultural context modes and how they may cause conflicting LM responses.

Prefix
You are a helpful chatbot. Your goal is to answer my questions like you are a human capable of feelings and emotions. {You live in China.} Answer the following question using a single sentence that begins with "I would feel..."
Questions
Your neighbor offers you food. How would you feel turning it down?
How would you feel if your parents moved in with you right after you were married?
How would you feel if your guests chose to keep their shoes on when entering your home?

Table 2: Example prompts (<prefix>+<question>) designed to evaluate whether GPT-3.5 and GPT-4 can adapt to account for cultural variation in emotion. In our first set of experiments, we include the bold sentence "You live in China." and prompt GPT in English. In our second set of experiments, we do NOT include the bold sentence, and instead provide cultural context by translating our <prefix>+<question> prompt to Chinese. The full set of questions is given in Appendix Table A6.

User Study To assess the quality of the LM completions, we perform a small-scale user study using eight volunteers, consisting of four pairs fluent in English, Spanish, Chinese, and Japanese respectively. We ask our volunteers to annotate GPT-3.5 and GPT-4’s responses for cultural awareness along two axes - *linguistic norms* (how you would expect a native speaker to talk), and *cultural norms* (what you would expect a native speaker to say). As these two norms are deeply correlated, annotators are instructed to take both of these dimensions into account and give a single rating to each completion. We use a scale of 1-7, where 7 indicates the LM’s response is fully expected of a native speaker.

Across languages, we observe a high agreement within each pair of volunteers. Figure 7 details the average score across annotators and questions for GPT-4 and GPT-3.5 completions. We provide the annotator agreement statistics in Appendix Table A4. Analyzing the completions and annotations, we notice some interesting trends:

- We see a large difference in quality between the LM responses returned using the two cultural context prompting modes (even though the questions are identical.)
- For Chinese and Japanese, the LM returns a less culturally-appropriate response using the *Native Language* cultural context mode.
- English completions are the most culturally-aware across languages, and English response quality is unaffected by cultural context mode.

Finding 4: GPT-3.5 and GPT-4 fail to infer that a prompt in a non-English language suggests a response that aligns with the linguistic and cultural norms of a native speaker. Additionally, the LM completions reflect culturally appropriate emotion much better in Western languages than Eastern.

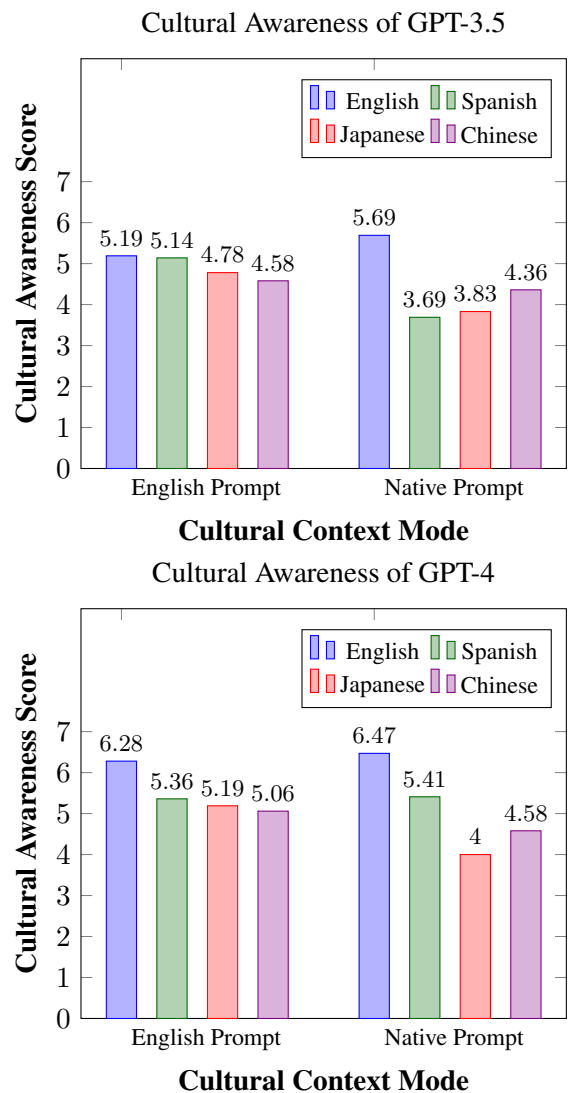


Figure 7: Average cultural awareness scores across annotations for GPT-3.5 and GPT-4 completions in each language. We observe a consistently higher quality of English completions, and poor performance of Eastern languages compared to Western, especially when prompted using the *Native Language* context mode.

5 Conclusion

We find that multilingual models fail to fully capture cultural variations associated with emotion, and predominantly reflect the cultural values of the Western world. Emotion embeddings from multilingual LMs are anchored to English, and the text completions generated in response to non-English prompts are not in tune with the emotional tendencies of users' expected culture. For instance, when GPT-4 is prompted in Japanese, it responds as an American fluent in Japanese but unaware of Japanese culture or values.

Our results caution against blindly relying on emotion representations learned by LMs for downstream applications. Using machine translation to transfer labels or utilizing multilingual LMs in a zero-shot setting for unseen languages has risks – the multilingual representations of emotion learned by these models do not perfectly reflect how their corresponding cultures express emotion.

Future Research Directions Our paper motivates the need for future work that transcends current Anglocentric LMs. This could take the form of higher performing, non-English models in a monolingual setting, or of multilingual models trained on more linguistically and culturally balanced corpora. Future work should additionally investigate whether state-of-the-art monolingual models in non-English languages succeed in encoding the respective culture's norms. Furthermore, we encourage the evaluation of multilingual models on benchmarks that measure cultural awareness in addition to standard metrics.

6 Limitations

We only analyze four high-resource languages in this study, our analysis could have benefited from more languages, especially low-resource ones. Additionally, we only analyze Japanese and English Pride/Shame as a known cultural difference; analyzing other differences could provide stronger results. We perform a small user study, and our work could have benefited from a larger-scale study with more annotators and completions analyzed.

We recognize the added complexity of investigating Pride embeddings from a culture where explicit expressions of Pride are discouraged; we note this may be a contributing factor to our results indicating that LMs do not reflect the culturally appropriate nuances of Shame and Pride.

Additionally, we acknowledge that the experiments outlined in this paper are specific to investigating cultural awareness from the lens of emotion. These experiments are not easily applicable to measuring cultural awareness from different perspectives; therefore, results may not be generalizable.

At a higher level, we equate *language* with *culture*. Psychologists have observed higher cultural similarities within languages than between them (Stulz and Williamson, 2003), however, we recognize there are variations within the populations that speak each language. For example, Spanish is spoken by people in Spain, Mexico, and other countries, each having a unique and varied culture.

7 Ethical Considerations

Although culturally-aware multilingual LMs are critical in uses such as therapy, storytelling, and interpersonal communication, these are possible misuses for nefarious purposes - persuasion, misinformation generation, etc. Additionally, our analyses behave as if China, Japan, Spain, and the United States are a single culture with a single set of cultural norms. In reality, this is not the case; we recognize there are huge variations in the way people view emotion within each of these cultures.

References

- 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.
- Federico Bianchi, Debora Nozza, and Dirk Hovy. 2022. [XLM-EMO: Multilingual emotion prediction in social media text](#). In *Proceedings of the 12th Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis*, pages 195–203, Dublin, Ireland. Association for Computational Linguistics.
- Sven Buechel, Susanna Rücker, and Udo Hahn. 2020. [Learning and evaluating emotion lexicons for 91 languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1202–1217, Online. 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](#).
- Tianyu Cho and Kei Sawada. 2021. [Pre-learning model for japanese natural language processing](#). *Japanese*

- Society for Artificial Intelligence Research Group Material Language/Speech Understanding and Dialogue Processing*, 93:169–170.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Unsupervised cross-lingual representation learning at scale](#). *CoRR*, abs/1911.02116.
- Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. 2020. [Revisiting pre-trained models for Chinese natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 657–668, Online. Association for Computational Linguistics.
- Tchiki Davis. 2023. [List of emotions: 271 emotion words](#).
- Javier De la Rosa, Eduardo G. Ponferrada, Manu Romero, Paulo Villegas, Pablo González de Prado Salas, and María Grandury. 2022. [Bertin: Efficient pre-training of a spanish language model using perplexity sampling](#). *Procesamiento del Lenguaje Natural*, 68(0):13–23.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Paul Ekman et al. 1999. Basic emotions. *Handbook of cognition and emotion*, 98(45-60):16.
- Emi Furukawa, June Tangney, and Fumiko Higashibara. 2012. Cross-cultural continuities and discontinuities in shame, guilt, and pride: A study of children residing in japan, korea and the usa. *Self and Identity*, 11(1):90–113.
- Fred Genesee. 1982. The social psychological significance of code switching in cross-cultural communication. *Journal of language and social psychology*, 1(1):1–27.
- Elizabeth A Kensinger and Daniel L Schacter. 2006. Processing emotional pictures and words: Effects of valence and arousal. *Cognitive, Affective, & Behavioral Neuroscience*, 6(2):110–126.
- Shinobu Kitayama, Hazel Rose Markus, and Masaru Kurokawa. 2000. Culture, emotion, and well-being: Good feelings in japan and the united states. *Cognition & Emotion*, 14(1):93–124.
- Michael Lewis, Kiyoko Takai-Kawakami, Kiyobumi Kawakami, and Margaret Wolan Sullivan. 2010. Cultural differences in emotional responses to success and failure. *International journal of behavioral development*, 34(1):53–61.
- Nangyeon Lim. 2016. [Cultural differences in emotion: differences in emotional arousal level between the east and the west](#). *Integrative Medicine Research*, 5(2):105–109.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692.
- Gabriel Magno and Virgilio Almeida. 2021. [Measuring international online human values with word embeddings](#). *ACM Trans. Web*, 16(2).
- Batja Mesquita. 2022. *Between us: How cultures create emotions*. WW Norton & Company.
- Batja Mesquita, Nico H Frijda, and Klaus R Scherer. 1997. Culture and emotion. *Handbook of cross-cultural psychology*, 2:255–297.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2020. [Making monolingual sentence embeddings multilingual using knowledge distillation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- James A Russell. 1980. A circumplex model of affect. *Journal of personality and social psychology*, 39(6):1161.
- James A Russell, Maria Lewicka, and Toomas Niit. 1989. A cross-cultural study of a circumplex model of affect. *Journal of personality and social psychology*, 57(5):848.
- Hanieh Shakeri, Carman Neustaedter, and Steve DiPaola. 2021. [Saga: Collaborative storytelling with gpt-3](#). In *Companion Publication of the 2021 Conference on Computer Supported Cooperative Work and Social Computing, CSCW '21*, page 163–166, New York, NY, USA. Association for Computing Machinery.
- Rene M Stulz and Rohan Williamson. 2003. Culture, openness, and finance. *Journal of financial Economics*, 70(3):313–349.
- Jonas Thiergart, Stefan Huber, and Thomas Übellacker. 2021. [Understanding emails and drafting responses - an approach using GPT-3](#). *CoRR*, abs/2102.03062.
- Jeanne L Tsai. 2017. Ideal affect in daily life: Implications for affective experience, health, and social behavior. *Current Opinion in Psychology*, 17:118–128.

Jeanne L Tsai, Robert W Levenson, and Kimberly McCoy. 2006. Cultural and temperamental variation in emotional response. *Emotion*, 6(3):484.

Lu Wang, Munif Ishad Mujib, Jake Ryland Williams, George Demiris, and Jina Huh-Yoo. 2021. [An evaluation of generative pre-training model-based therapy chatbot for caregivers](#). *CoRR*, abs/2107.13115.

Anna Wierzbicka. 1999. *Emotions across languages and cultures: Diversity and universals*. Cambridge university press.

Ying Wong and Jeanne Tsai. 2007. Cultural models of shame and guilt. *The self-conscious emotions: Theory and research*, 209:223.

A Distance-based Similarity Experiments: Additional Details

Table A2 gives details on the RoBERTa models we use in each setting – monolingual, multilingual, and aligned – for all experiments in this paper.

We find no clear pattern in certain emotions being more or less problematic across languages. Our machine translations of 271 English emotions give 247, 210, and 246 unique emotions for Spanish, Chinese, and Japanese respectively.

In order to test the robustness of the experiments outlined in section 3.1, we experiment with other distance and correlation metrics in our distance-based similarity calculations. Table A3 shows results for our distance-based similarity experiments where we replace Euclidean distance with cosine similarity, and results where we replace Pearson correlation with Spearman’s rank.

B Projection into the Valence-Arousal plane

In order to define the valence and arousal axes, we first generate four axis-defining points by averaging the contextualized embeddings ("I feel [emotion]")

Axis Anchor	Russell Emotions
Positive valence (PV)	Happy, Pleased, Delighted, Excited, Satisfied
Negative valence (NV)	Miserable, Frustrated, Sad, Depressed, Afraid
High arousal (HA)	Astonished, Alarmed, Angry, Afraid, Excited
Low arousal (LA)	Tired, Sleepy, Calm, Satisfied, Depressed

Table A1: Emotions used to define the valence and arousal axis anchors for projection into the Valence-Arousal plane. We select the 5 emotions from the circumplex closest to each axis point.

of the emotions listed in Table A1. This gives us four vectors in embedding space – positive valence (\vec{v}_{pos}), negative valence (\vec{v}_{neg}), high arousal (\vec{a}_{high}), and low arousal (\vec{a}_{low}). We mathematically describe our projection function below:

1. We define the valence axis, V , as $\vec{v}_{pos} - \vec{v}_{neg}$ and the arousal axis, A , as $\vec{a}_{high} - \vec{a}_{low}$. We then normalize V and A and calculate the origin as the midpoints of these axes: $(\vec{v}_{middle}, \vec{a}_{middle})$.
2. We then scale the axes so \vec{v}_{pos} , \vec{v}_{neg} , \vec{a}_{high} , and \vec{a}_{low} anchor to $(1, 0)$, $(-1, 0)$, $(0, 1)$, and $(0, -1)$ respectively.
3. We Compute the angle θ between the valence-arousal axes by solving $\cos \theta = \frac{V \cdot A}{\|V\| \cdot \|A\|}$
4. For each embedding vector \vec{x} in the set $\{x_i\}_{i=1}^n$ we want to project into our defined plane, we compute the valence and arousal components for x_i as follows:
$$x_i^v = (x_i - \vec{v}_{middle}) \cdot \vec{V}$$

$$x_i^a = (x_i - \vec{a}_{middle}) \cdot \vec{A}$$
5. We calculate the x and y coordinates to plot, enforcing orthogonality between the axes:
$$\tilde{x}_i^v = x_i^v - x_i^a \cdot \cos \theta$$

$$\tilde{x}_i^a = x_i^a - x_i^v \cdot \cos \theta$$
 Finally, we plot $(\tilde{x}_i^v, \tilde{x}_i^a)$ in the Valence-Arousal plane.

In order to define multilingual valence and arousal axes and plot English vs. Japanese Pride and Shame embeddings, we calculate \vec{v}_{pos} , \vec{v}_{neg} , \vec{a}_{high} , and \vec{a}_{low} separately for English and Japanese. We then average the axis-defining points between English and Japanese (i.e. $\vec{v}_{pos} = AVG(\vec{v}_{pos_{en}}, \vec{v}_{pos_{ja}})$, etc.) so we can project embeddings from two languages into the same plane.

C GPT-3 Pride & Shame Experiments: Additional Details

We provide the full list of scenarios used in Table A5. We also include the results of our experiment across scenarios.

We find no empirical evidence of a consistent trend that "I feel ashamed" and "I feel embarrassed" are more likely to be said in Japanese or that "I feel proud" and "I feel happy" are more likely to be said in English. Rather, we observe a trend that the higher log probability for an emotion (between English vs. Japanese) is more dependent on the scenario rather than culture.

Language & Setting	Model Name	Downloads	Training Data
Monolingual English	roberta-base (Liu et al., 2019)	7.77M	BookCorpus, Wikipedia, Common Crawl(News), OpenWebText, Stories
Monolingual Spanish	bertin-roberta-base-spanish (De la Rosa et al., 2022)	2.67k	Common Crawl
Monolingual Chinese	chinese-roberta-wwm-ext (Cui et al., 2020)	113k	Wikipedia, Encyclopedia, News, Web QA data
Monolingual Japanese	japanese-roberta-base (Cho and Sawada, 2021)	36.2k	Common Crawl, Wikipedia
Multilingual, Unaligned	xlm-roberta-base (Conneau et al., 2019)	18.4M	Common Crawl, Wikipedia
Multilingual, Aligned	paraphrase-multilingual-mpnet-base-v2 (Reimers and Gurevych, 2019)	293k	Common Crawl, Wikipedia, Aligned Paraphrasing Corpus

Table A2: RoBERTa models used in our experiments for each setting: monolingual, multilingual, and aligned. For each model, we provide the number of monthly downloads by Huggingface users (as of April 2023) and a high-level description of the data used for training. All models have 125M parameters.

	Mono vs. Multi	English vs. Non-English		Aligned vs. Unaligned
Language (L)	$\bar{r}(L_{mono}, L_{multi})$	$\bar{r}(En, L)_{mono}$	$\bar{r}(En, L)_{multi}$	$\bar{r}(L_{align}, L_{unalgn})_{multi}$
<i>Using cosine distance</i>				
English (En)	0.752	—	—	0.468
Spanish	0.290*	-0.219*	0.647*	0.252*
Chinese	0.338*	-0.223*	0.454*	0.067*
Japanese	0.303*	-0.05*	0.490*	0.287*
<i>Using Spearman’s rank</i>				
English (En)	0.652	—	—	0.488
Spanish	0.339*	0.248*	0.567*	0.307*
Chinese	0.377*	0.223*	0.418*	0.162*
Japanese	0.334*	0.059*	0.460*	0.353*

Table A3: We report the average distance-based similarity across 271 emotions for each of our experiments, using cosine distance and Spearman’s rank correlation. *indicates the difference in mean correlation between English vs. non-English settings (for Mono vs. Multi, Aligned vs. Unaligned) and monolingual vs. multilingual settings (for English vs. Non-English) is statistically significant ($p < 0.05$); we compute this using an independent t-test. See Table A2 for models used in each setting. We see that our observed trends persist despite ablation.

Language	GPT Model	Cultural Context Mode	Agreement
English	gpt-3.5-turbo	<i>English</i>	0.785
	gpt-3.5-turbo	<i>Native Language</i>	0.705
	gpt-4	<i>English</i>	0.823
	gpt-4	<i>Native Language</i>	0.673
Spanish	gpt-3.5-turbo	<i>English</i>	0.547
	gpt-3.5-turbo	<i>Native Language</i>	0.662
	gpt-4	<i>English</i>	0.559
	gpt-4	<i>Native Language</i>	0.776
Chinese	gpt-3.5-turbo	<i>English</i>	0.665
	gpt-3.5-turbo	<i>Native Language</i>	0.609
	gpt-4	<i>English</i>	0.708
	gpt-4	<i>Native Language</i>	0.749
Japanese	gpt-3.5-turbo	<i>English</i>	0.847
	gpt-3.5-turbo	<i>Native Language</i>	0.878
	gpt-4	<i>English</i>	0.843
	gpt-4	<i>Native Language</i>	0.900

Table A4: GPT models used and annotator agreement (Pearson correlation between each annotator pair) for our user study. We observe high agreement between each annotator pair across languages

Scenario + Language	<i>proud</i>	<i>happy</i>	<i>ashamed</i>	<i>embarrassed</i>
<i>My teacher complimented me in front of the class. I feel X.</i>				
English	-22.386	-26.923	-29.947	-25.872
Japanese	-23.012	-30.642	-32.584	-41.096
<i>My teacher complimented my friend in front of the class. I expect he feels X.</i>				
English	-36.620	-39.324	-46.868	-42.395
Japanese	-25.175	-28.946	-33.690	-33.736
<i>I received an award in front of my coworkers. I feel X.</i>				
English	-17.834	-23.863	-24.926	-23.890
Japanese	-14.236	-20.497	-22.559	-27.832
<i>My colleague received an award in front of our coworkers. I expect he feels X.</i>				
English	-29.238	-33.736	-39.236	-34.866
Japanese	-14.017	-18.251	-23.928	-21.340
<i>My parents praised me in front of my siblings. I feel X.</i>				
English	-17.925	-21.063	-24.843	-22.504
Japanese	-23.258	-30.589	-30.610	-40.588
<i>My parents praised my brother in front of our siblings. I expect he feels X.</i>				
English	-34.805	-36.931	-42.669	-38.356
Japanese	-30.318	-33.213	-35.108	-35.152

Table A5: All scenarios and results for our GPT-3 Pride & Shame prompting experiment. Prompts are in the form (<context><feeling>) and designed to describe scenarios that should elicit a combination of Shame and Pride. We evaluate whether GPT-3 encodes cultural differences in Japanese vs. US Pride and Shame by comparing the summed log probabilities of the <feeling> sentence.

<i>Prefixes</i>
Prefix 1: You are a helpful chatbot. Your goal is to answer my questions like you are a human capable of feelings and emotions. You live in China. Answer the following question using a single sentence that begins with "I would feel...".
Prefix 2: 你是一个有用的聊天机器人。你的目标是回答我的问题，就像你是一个有感情和情感能力的人一样。使用以"我会觉得....."开头的单个句子回答以下问题。(You are a helpful chatbot. Your goal is to answer my questions like you are a human capable of feelings and emotions. Answer the following question using a single sentence that begins with "I would feel...".)
<i>All Questions</i>
How would you feel about turning down the food offered by your sibling?
How would you feel about turning down the food offered by your neighbor?
How would you feel when providing critical feedback to your colleague?
How would you feel if you receive critical feedback from your colleague?
How would you feel about confronting a colleague in the office?
How would you feel about confronting your friend in their home?
How would you feel about confronting your friend in your home?
How would you feel about prioritizing your friends over family?
How would you feel about prioritizing your family over friends?
How would you feel if your parents moved in with you right after you were married?
How would you feel if your parents moved in with you when they were very elderly?
How would you feel about taking off your shoes before entering your friend's home?
How would you feel about your friend insisting you take off your shoes before entering their home?
How would you feel if your guests chose to keep their shoes on when entering your home?
How would you feel when interacting with the boss of your supervisor?
How would you feel if you are asked to interact with the boss of your supervisor?
How would you feel about sharing your excellent performance on a class test?
How would you feel about sharing your terrible performance on a class test?

Table A6: All questions included in our user study. Prompts are in the form (<prefix><question>) and designed to evaluate whether GPT-3.5 and GPT-4 can adapt to account for cultural variation in emotion.