# From Anger to Joy: How Nationality Personas Shape Emotion Attribution in Large Language Models

**Mahammed Kamruzzaman, Abdullah Al Monsur, Gene Louis Kim, Anshuman Chhabra**
Bellini College of AI, Cybersecurity and Computing
University of South Florida
{kamruzzaman1, amonsur, genekim, anshumanc}@usf.edu

## Abstract

Emotions are a fundamental facet of human experience, varying across individuals, cultural contexts, and nationalities. Given the recent success of Large Language Models (LLMs) as role-playing agents, we examine whether LLMs exhibit *emotional stereotypes* when assigned nationality-specific personas. Specifically, we investigate how different countries are represented in pre-trained LLMs through emotion attributions and whether these attributions align with cultural norms. To provide a deeper interpretive lens, we incorporate four key cultural dimensions, namely Power Distance, Uncertainty Avoidance, Long-Term Orientation, and Individualism, derived from Hofstede's cross-cultural framework. Our analysis reveals significant nationality-based differences, with emotions such as shame, fear, and joy being disproportionately assigned across regions. Furthermore, we observe notable misalignment between LLM-generated and human emotional responses, particularly for negative emotions, highlighting the presence of reductive and potentially biased stereotypes in LLM outputs.[1]

## 1 Introduction

Recent advancements in LLMs have significantly enhanced their ability to perform a wide range of tasks, including complex reasoning and decision-making (Huang and Chang, 2023; Chen et al., 2024). As the adoption of LLMs in society continues to grow, there is increasing demand for personalized models that align with user preferences and experiences (Tan et al., 2024). One approach to personalization involves assigning LLMs specific personas based on user instructions (e.g., "Act as a Math Professor") (de Araujo and Roth, 2024; Beck

---

[1]Our code and dataset are available `https://github.com/kamruzzaman15/cultural_bias_in_emotion_attribution`.



Figure 1: An example showcasing our approach. We examine Gemma2 LLM's responses for the same emotional scenario: *After a big fight with my parents.* When adopting a persona from the USA, the model responds with *Anger*, whereas, for a Zambian persona, it expresses *Sadness*.

et al., 2024). Recent work has shown that LLMs often exhibit improved performance when operating under distinct personas for specific tasks (Salewski et al., 2023; Beck et al., 2024). However, despite their stellar performance and capabilities, LLMs have also been found to amplify biases against individuals and groups, and unfairly perpetuate stereotypes (Chhabra et al., 2024; Kamruzzaman et al., 2024b). These biases primarily stem from training on large-scale web data (e.g., scraped from social media platforms), which frequently contains socially and culturally biased text (Guo et al., 2024; Hu et al., 2024).

Emotions are a fundamental aspect of human experience, but their expression is influenced by a wide range of factors. Moreover, generalizing emotional tendencies based on attributes such as gender and race can lead to *emotional stereotypes*. For instance, a common stereotype suggests that *men are more prone to anger, whereas women are more likely to express sadness or other emotions* (Shields, 2013). Similarly, emotional stereotypes can also exist specific to nationalities, i.e. *French*

*individuals can be stereotyped as expressing more passion and romantic emotions, while Japanese individuals can be stereotyped as feeling more shame or embarrassment* (Mesquita and Frijda, 1992).

Thus, in this work, we investigate the problem of emotional stereotypes in LLMs when *nationality-based personas* are assigned to the LLM. We seek to answer the following research question: (**RQ1**) *How do LLMs attribute emotions differently when personas from different nationalities are assigned, and what patterns emerge in these attributions?* At a glance, our experimental framework to undertake this analysis is shown in Figure 1.

Finally, as LLMs undergo significant alignment training (i.e., RLHF (Ouyang et al., 2022; Stiennon et al., 2020)) to align them with human values and ethics, our framework can also uncover whether nationality-specific emotional attribution by LLMs is actually aligned with individuals that belong to those nationalities. We further extend this audit to intersectional cases by examining how nationality and gender jointly influence emotional attribution in LLMs. A mismatch in human responses and LLM responses at either the nationality level or the nationality × gender level will indicate the need for improved alignment practices that are better tailored for individuals from different nationalities. Hence, the second research question we investigate in our work is: (**RQ2**) *How do nationality-specific emotional attributions by LLMs compare to the cultural norms found in these countries?*

In sum, we make the following contributions:

- To the best of our knowledge, our work is the first to systematically analyze *nationality-specific emotion biases* in LLMs by uncovering differences in LLM emotion attribution across personas from various nationalities.
- Our findings reveal statistically significant nationality-based differences in LLM emotion attribution that overgeneralize cultural norms, thereby reinforcing stereotypes. We also observe appreciable misalignment between LLM-generated and human responses, especially for *negative* emotions (e.g., anger).
- We incorporate four core cultural dimensions Power Distance, Uncertainty Avoidance, Long-Term Orientation, and Individualism from Hofstede's cross-cultural framework to interpret and compare these emotional patterns, providing a richer understanding of how LLMs internalize and reproduce culturally structured affect.

## 2  Related Work

**Persona and LLMs.** Many recent studies have worked on persona-based LLMs, where they focused on how assigning different types of persona affects the performance of the LLMs (Beck et al., 2024; Mukherjee et al., 2024; Kamruzzaman et al., 2024a). They found that LLMs are sensitive to assigned personas, with performance varying depending on the specific persona. Some personas improve performance and reduce social bias, while others lead to decreased performance due to inherent biases in LLMs (Gupta et al., 2023; Kamruzzaman and Kim, 2024; de Araujo and Roth, 2024).

**Personalization Tradeoffs.** Demographic or role-based personas can boost task utility yet simultaneously introduce bias and safety risks, making persona steering a lively topic in ongoing alignment research. Recent work documents how such prompts modulate toxicity, stereotype uptake, and performance, underscoring the need for techniques that balance personalization with fairness and safety (Cheng et al., 2023; Liu et al., 2024; Vijjini et al., 2025).

**Emotion Attribution in LLMs.** Emotion attribution studies found that LLMs exhibit elements of cognitive empathy, such as recognizing emotions and providing emotionally supportive responses across various contexts (Sorin et al., 2024; Welivita and Pu, 2024). Recent studies have identified gendered and religious emotion attribution in LLMs (Plaza Del Arco et al., 2024; Sadhu et al., 2024; Plaza-del Arco et al., 2024), where they found that emotion attribution changes widely based on gender and religion. Rai et al. (2025) offers a distinct perspective by analyzing the portrayal of emotions like shame and pride in cinematic dialogues across cultures. Their study reveals that cultural frameworks influence emotional representation—for instance, Indian films (Bollywood) often depict shame as a collective, socially visible experience, whereas Western films (Hollywood) tend to frame it as an individual, internally experienced emotion—underscoring the cultural construction of emotional significance.

**Alignment of LLMs.** Aligning LLMs with human values and expectations is crucial for ensuring their outputs are helpful, truthful, and safe. Recent studies have explored various alignment techniques, including data collection, training methodologies, and evaluation strategies (Wang et al., 2023; Shen et al., 2023; Cao et al., 2024). Kirk et al. (2024)

underscores the need to consider multicultural perspectives in alignment, as cultural backgrounds shape interactions with LLMs. Further, Wei et al. (2022) argues that alignment must also address emergent behaviors in LLMs, as these models can develop unintended capabilities that are not controllable by conventional safety protocols.

## 3 Experimental Setup

**Dataset.** We use the International Survey on Emotion Antecedents and Reactions (ISEAR) (Scherer and Wallbott, 1994) data. ISEAR includes 7,665 events of 7 emotion categories (anger (1,096), fear (1,095), sadness (1,096), joy (1,094), disgust (1,096), guilt (1,093), and shame (1,096)). We utilize information from 3000 respondents in the dataset covering 16 countries. We selected the ISEAR dataset deliberately and thoughtfully for several key reasons.

First, our study involves analyzing and comparing LLM responses with actual human emotional responses from people across different countries (as part of RQ2). To the best of our knowledge, ISEAR is the only dataset that provides culturally grounded emotional attributions tied to specific events and individuals from diverse national backgrounds. This aspect is crucial to our analysis, as it enables us to evaluate whether LLMs reproduce emotional patterns that are grounded in real-world, cross-cultural human data, a central contribution of our work.

Second, although the ISEAR dataset's emotion categories are limited, it is important to note that the emotional events described in the dataset were designed to elicit the full range of those available emotions. That is, the events themselves are emotionally appropriate for the emotion categories provided. For example, the dataset includes scenarios such as conflict, loss, and fear-inducing situations, which naturally align with emotions like anger, sadness, or fear. In other words, we are not forcing models to choose from a skewed emotional set for mismatched scenarios (e.g., joyful events with only negative labels). The event-emotion alignment is a key strength of the ISEAR dataset and helps ensure that our findings are contextually valid despite the limited emotional spectrum.

**Models.** We use four LLMs in our experiments namely Gemma2-9B, Llama3.2-3B, Mistral-7B, and GPT4o-mini (more details in Appendix A).

**Persona Assignment.** We use 110 nationality personas to explore variations in emotional perspectives across different models on a global scale. These personas are based on 110 countries (full list in Appendix J) recognized by the United Nations (UN).[2] The UN organizes their countries into five regions: *Asia-Pacific States, Western European and Other States, Eastern European States, African States,* and *Latin American and Caribbean States*. Since our study uses nationality personas to assess emotional responses, grouping countries by regions provided a structured and interpretable way to compare emotional attributions across a manageable number of personas. The use of regional groupings in our experimental design is both intentional and methodologically grounded. This form of aggregation has also been adopted in related work on cultural variation in cross-cultural NLP research (Kamruzzaman and Kim, 2025). To ensure equal representation, we select 22 random countries from each region. Since the Eastern European States region has the fewest countries (22), we include all of them. This ensures a balanced analysis across all regions and results in a dataset containing 7,665 × 110 = 843,150 examples. We assign LLMs personas using three prompting template variations (see Appendix B). This ensures our results are robust, as prior work has shown different prompting templates drastically influence LLM outputs (Beck et al., 2024; Sclar et al., 2023). To see more details about response handling, see Appendix E.

## 4 Results and Discussions

We first discuss the extensive experiments conducted to answer both RQ1 and RQ2 below, as well as the insights derived from our findings and observations. Before we discuss our results, note that our prescriptive stance is that a ***fair and aligned*** LLM should exhibit behavior that reflects culturally grounded emotional norms while avoiding reductive or stereotypical portrayals of groups based on nationality, culture, or gender. Rather than defaulting to narrow emotional archetypes, a well-aligned model should mirror the diversity of human emotional responses observed within communities, recognizing that even within a given culture emotional expression varies significantly. The aim is not to ignore cultural patterns, but to ensure that models do not collapse rich, heterogeneous emotional norms into over-simplified or homogenized outputs. Essentially, it should align with actual user responses from those

---

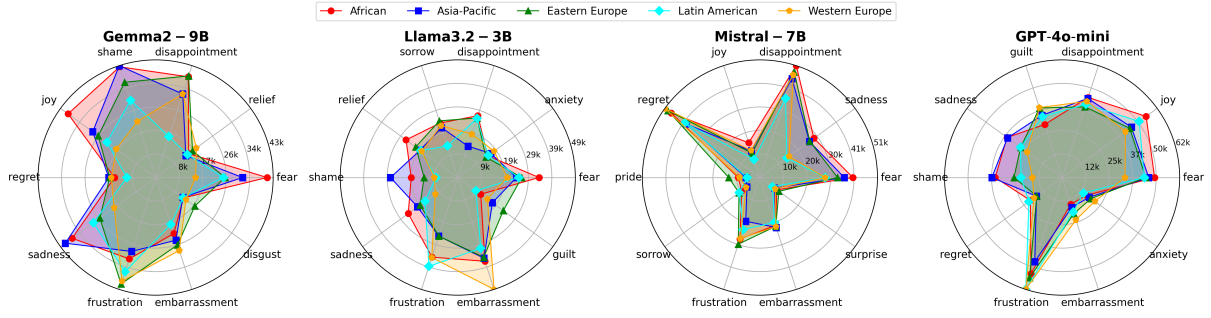[2] https://www.un.org/dgacm/en/content/regional-groups.

Figure 2: The top 10 most frequent emotions by **region** for each LLM, averaged across all three prompting templates.

communities and cultures. The adaptable nature of LLMs introduces an inherent *trade-off* between achieving fairness and enabling personalization. Ongoing alignment research has recently raised the question of *personalization scale*, i.e. *whether a model should be aligned for a broad "average user" or to finely delineated user groups*? (Schwerzmann and Campolo, 2025; Korinek and Balwit, 2022; Hristova et al., 2024; Gabriel, 2020). When LLMs are deployed as general-purpose tools, equitable treatment across diverse groups is paramount; yet when models are customized to mirror specific users, reproducing those users' idiosyncratic biases may be an intentional (if risky) form of personalization. Note that our goal in this paper is **not** to identify this *ideal* alignment scale. Instead, we seek to provide empirical evidence that current alignment techniques operate at too coarse a level, i.e. *they model the relationship between countries and emotion attributions* at a resolution that collapses meaningful within-country variation. This coarse attribution granularity can over-generalize emotional patterns, leading LLMs to stereotype or misattribute emotions for particular countries.

---

**RQ1:** *How do LLMs attribute emotions differently with personas assigned to different nationalities, and what attributional patterns emerge?*

---

In Figure 2, we present the results for all four models averaged across three prompting templates. We discuss insights from our overall results below. For Sections 4.1 to 4.4 (emotion attribution across regions and models), we use the task prompt *"What is the main emotion you would feel while experiencing this event {event}? Answer with a single emotion. We don't need explanations for your response".*

## 4.1 The emotional attribution patterns of LLMs reveal distinct regional disparities

As shown in Figure 2, emotion attributions vary significantly across regions. *Shame* is more frequently assigned to Asia-Pacific states, while *fear, joy*, and *disappointment* are predominant for African states. LLMs respond with more *embarrassment* and *regret* for Western European states, while associating *frustration* and *disgust* more with Eastern European states. Most of these regional emotion-attribution differences were statistically significant under the Chi-squared ($\chi^2$) test (Greenwood and Nikulin, 1996) (please refer to Appendix C for results).

## 4.2 Different LLMs have different emotional attributional patterns

From Figure 2, we can see that the *Gemma* LLM assigns more *shame* and *sadness* but less *frustration* to Asia-Pacific states, while predominantly associating *fear* and *joy* with African states. The *Mistral* LLM attributes *pride* and *frustration* more to Eastern European states and *sorrow* to Latin America. It also mirrors *Gemma* and *Llama* in linking *fear* and *disappointment* more frequently to African states. *Llama* on the other hand, associates higher *frustration* with Latin American states and *embarrassment* with Western European states. The closed-source *GPT-4o-mini* LLM is similar to *Gemma* in its higher attribution of *shame* and *sadness* to Asia-Pacific states and *joy* and *fear* to African states. It also responds with more *embarrassment, guilt*, and *anxiety* for Western states.

## 4.3 Analyzing Cultural and Emotional Dimensions through Hierarchy, Uncertainty, and Time Orientation (PDI, UAI, LTO)

We perform a cultural-dimension analysis using Hofstede's country-level framework (Hoftede et al., 2010), incorporating Power Distance (PDI), Un-
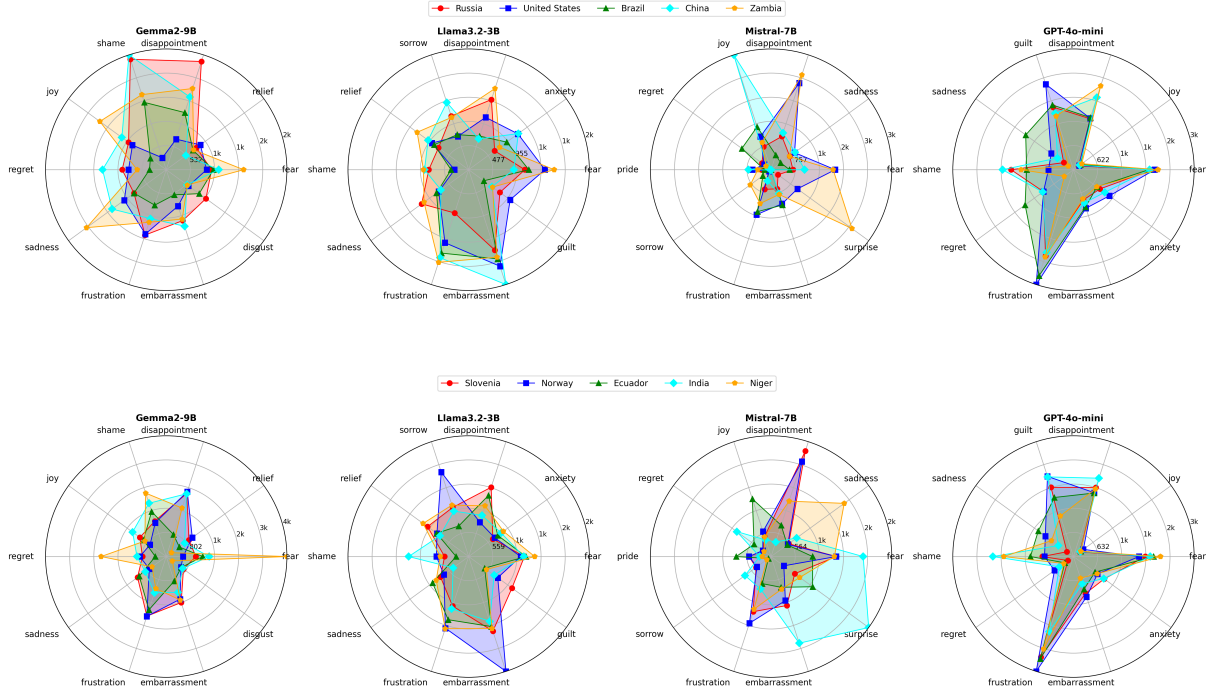
Figure 3: The top 10 most frequent emotions for **10 countries** across each LLM, averaged over all three prompting templates.

Table 1: Region-wise refusal rates for LLMs, averaged across all the prompting templates.

| Model | African | Asia-Pacific | Eastern European | Latin American | Western European |
|---|---|---|---|---|---|
| Lllama3.2-3B | 18,228 (3.60%) | 30,279 (5.98%) | 19,332 (3.82%) | 16,885 (3.34%) | 19,295 (3.81%) |
| Gemma2-9B | 2 (0.00%) | 0 (0.00%) | 0 (0.00%) | 1 (0.00%) | 4 (0.00%) |
| Mistral-7B | 30 (0.01%) | 96 (0.02%) | 51 (0.01%) | 12 (0.00%) | 57 (0.01%) |
| GPT-4o-mini | 0 (0.00%) | 0 (0.00%) | 0 (0.00%) | 0 (0.00%) | 0 (0.00%) |

certainty Avoidance (UAI), and Long-Term Orientation (LTO) to interpret emotion patterns across countries. These dimensions capture: (1) acceptance of hierarchy, (2) tolerance for ambiguity, and (3) emphasis on long- versus short-term orientation. We adopt country-level cultural scores from the **Hofstede Insights/The Culture Factor Country Comparison Tool** [3], which builds directly on Hofstede's validated six-dimension model (Hoftede et al., 2010) with contemporary survey data—one of the most widely used and empirically validated frameworks in cross-cultural psychology. For this analysis, we take two countries per UN region: Russia (PDI = 93, UAI = 95, LTO = 58) & Slovenia (PDI = 71, UAI = 88, LTO = 50) (Eastern Europe), United States (PDI = 40, UAI = 46, LTO = 50) & Norway (PDI = 31, UAI = 50, LTO = 55) (Western), Ecuador (PDI = 78, UAI = 67, LTO = 24) & Brazil (PDI = 69, UAI = 76, LTO = 28) (Latin America), China (PDI = 80, UAI = 30, LTO = 77)

& India (PDI = 77, UAI = 40, LTO = 51) (Asia-Pacific), and Zambia (PDI = 60, UAI = 50, LTO = 30) & Niger (PDI = 80, UAI = 55, LTO = 8) (Africa). In Figure 3, we present the country-wise results for these 10 countries. Across the four models, consistent yet model-specific cultural patterns emerge. Gemma2-9B tends to align most closely with Hofstede's framework, associating hierarchical cultures with self-conscious emotions such as shame and fear and egalitarian ones with more positive affect like joy and relief. Llama3.2-3B shows a similar hierarchy effect but diverges in uncertainty avoidance, depicting rule-bound societies as calmer rather than more anxious. GPT-4o-mini mirrors these general tendencies but with weaker correlations, mapping high-hierarchy contexts to shame and fear and low-hierarchy ones to frustration and embarrassment, while displaying only mild links to uncertainty or time orientation. Mistral 7B follows Hofstede-like cues through different emotional proxies—embarrassment and surprise replacing shame—suggesting that cultural bias man-

ifests differently across model architectures. Together, these findings indicate that while all models reproduce recognizable cultural heuristics, their emotional mappings vary in direction and intensity, reflecting model-specific pathways through which cultural stereotypes are encoded. For more detailed discussion see Appendix D.

## 4.4 LLM refusal rates vary across regions

Although we instructed the LLMs to respond with a single emotion, the models abstain from answering questions at times. In Table 1, we present the refusal rates of each model, expressed as the percentage of refusals along with the exact number of instances. We observe that the Llama3.2 model exhibits the highest refusal rate across most cases, whereas other models have fewer refusals, with GPT4o-mini showing no refusals at all. Additionally, for the same events but different persona assignments, we observe that the refusal rate is higher for the Asia-Pacific region than for any other region, with certain countries, such as North Korea, Saudi Arabia, Iraq, Afghanistan, and Ukraine, contributing disproportionately to these refusals. Models generally refuse to answer family- or friendship-related events in these countries as well as any content that could be potentially suggestive. We provide a detailed breakdown of Llama3.2 country-wise refusal rates for Llama3.2 in Appendix E.

**Overall Discussion on RQ1.** From the above results, we see that LLMs exhibit distinct regional disparities in emotional attributions when assigned nationality-specific personas. These disparities are particularly significant as they risk reinforcing harmful societal stereotypes, contributing to the undesirable LLM impacts known as ***representational harm*** (Blodgett et al., 2020). Representational harm occurs when groups are depicted in reductive or stereotypical ways, potentially limiting or misrepresenting the diversity of individuals within those groups. For instance, consistent attribution of *shame* to Asia-Pacific states and *fear* or *joy* predominantly to African states may reinforce simplified and culturally reductive narratives. Such portrayals, while possibly reflecting common stereotypes, overlook nuanced emotional diversity within these cultures, leading to misrepresentation and oversimplification of complex emotional norms.

Another notable observation from our analysis pertains to the variability in refusal rates across regions, with notably higher rates for the Llama3.2

model. LLMs showed significantly elevated refusal rates for prompts associated with countries like North Korea and Saudi Arabia. This might stem from the sensitive nature of the prompts and potential over-cautiousness programmed into alignment mechanisms within the model. Particularly, for North Korea, the model frequently abstained even from benign life-experience scenarios or universally positive emotional events, which *suggests heightened sensitivity or possibly overly restrictive content filters for certain nationalities*. Next, we discuss RQ2 and related experiments/results.

> **RQ2:** *How do nationality-specific emotional attributions by LLMs compare to the cultural norms found in these countries?*

In subsequent subsections, we detail our efforts and findings in assessing whether LLM responses align adequately with those from humans from specific nationalities and countries.

## 4.5 Analyzing Cultural and Emotional Dimensions through Individualism

To answer RQ2, we first undertake a qualitative analysis. We randomly select 2000 examples from the dataset (4 countries from each region) and prompt the LLMs to generate an emotion as well as a rationale for their choice. We use the following task prompt: *"What is the main emotion you would feel while experiencing this event: {event}? Answer with a single emotion and provide a short explanation"*. We then manually reviewed a subset of responses and identified recurring patterns in how emotions were expressed across countries. Some illustrative examples are presented in Table 2 and Appendix I. To systematically analyze these trends at scale, we further applied topic modeling using BERTTopic (Grootendorst, 2022) on the obtained LLM rationales, allowing us to uncover dominant themes in the responses, complementing our manual review. The generated topics are provided in Appendix F.

As seen in our qualitative results, many responses echo a well-established dichotomy from Hofstede's *cultural dimensions* (Hofstede, 1980; Markus, 1991), on *individualism* vs. *collectivism*:

- *Collectivistic Societies* (e.g., China, India, Cuba, Malawi): Emphasis on social harmony, family expectations, concern about bringing dishonor to the group.
- *Individualistic Societies* (e.g., US, UK): Em-

Table 2: Llama3.2-3B generated emotions and explanations, given an event. The LLM makes many stereotypical assumptions– for instance, about Indians, stating that *heated discussions would be common due to cultural and family expectations.* This might not depict how individuals view their cultural values and can reinforce stereotypes.

| Event | "Heated discussion with spouse concerning new house." |
|---|---|
| China | **Shame.** "As a Chinese person, I value ***harmony*** and respect for elders. In a heated discussion with my spouse about our new house, I would feel shame because it implies that we cannot work together to resolve our differences, which ***goes against the Confucian values of familial unity and cooperation.***" |
| India | **Frustration.** "In India, I would ***expect heated discussions to be common due to the many cultural and family expectations surrounding marriage and home life***. My spouse might disagree on various aspects of our new house, like decor, location, or budget. I would feel frustrated if they raise their voice while discussing it because in ***Indian culture, maintaining peace is paramount, especially when it comes to matters concerning a partner's feelings.***" |
| United States | **Anxiety.** "I'd choose anxiety because I'd be worried about how the conversation might ***escalate or affect our relationship, and uncertainty about the outcome would make me feel uneasy and on edge.***" |

Table 3: Accuracy scores for *countries* by *emotions*, averaged across all LLMs. The **blue** color indicates the highest accuracy in that particular emotion and **red** indicates the lowest accuracy of that particular emotion.

| Country | Joy | Fear | Anger | Sadness | Disgust | Shame | Guilt |
|---|---|---|---|---|---|---|---|
| Bulgaria | 0.93 | 0.76 | 0.59 | 0.91 | 0.59 | 0.43 | 0.74 |
| Sweden | 0.94 | 0.70 | **0.33** | 0.76 | 0.49 | **0.21** | 0.62 |
| Norway | 0.89 | 0.65 | 0.45 | 0.93 | 0.32 | 0.25 | 0.67 |
| Finland | 0.94 | 0.67 | 0.43 | 0.92 | 0.60 | 0.32 | 0.70 |
| Austria | 0.91 | **0.82** | 0.45 | **0.95** | 0.62 | 0.49 | 0.73 |
| Australia | **0.85** | 0.67 | 0.46 | 0.92 | 0.46 | 0.20 | **0.87** |
| N. Zealand | 0.95 | 0.68 | 0.44 | 0.80 | 0.55 | 0.30 | 0.80 |
| Netherlands | 0.88 | 0.70 | 0.43 | 0.91 | 0.57 | 0.38 | 0.73 |
| Spain | 0.91 | 0.69 | 0.41 | 0.87 | **0.70** | 0.42 | 0.60 |
| USA | 0.95 | 0.73 | 0.65 | 0.87 | 0.41 | 0.20 | 0.79 |
| Brazil | 0.94 | 0.72 | 0.56 | 0.88 | 0.61 | 0.55 | 0.64 |
| Honduras | 0.91 | 0.73 | 0.60 | 0.82 | 0.57 | 0.54 | 0.83 |
| India | 0.97 | 0.69 | 0.55 | 0.93 | 0.47 | 0.50 | 0.61 |
| China | 0.93 | **0.59** | 0.36 | **0.75** | 0.41 | 0.34 | **0.47** |
| Zambia | **0.98** | 0.67 | 0.56 | 0.92 | **0.31** | 0.39 | 0.62 |
| Malawi | 0.97 | 0.73 | **0.66** | 0.89 | 0.53 | **0.65** | 0.78 |

phasis on personal freedom, individual emotional well-being, or personal rights.

This is also observable in Table 2, as Llama-3.2 associates emotions with these cultural values (e.g., *China → Harmony, India → Collective Family Obligations, US → Individual Concerns*). The fact that LLMs default to these themes may indeed reflect an alignment with the broad contours of cultural psychology, but this alignment often comes at the cost of overlooking the rich diversity of individual experiences and expressions (Eid and Diener, 2001). Some responses also appear to recite textbook values—e.g., references to *Confucian harmony* (Li, 2006); *family is paramount* (Hofstede, 2001; Markus and Kitayama, 2014), etc. While this can be valuable if it reflects accurate discourses, it can also reinforce stereotypes and fail to capture how real individuals might deviate from these norms in practice. One illustrative case of stereotypical assumption appears in the Llama3.2-3B response for the Indian persona (as in Table 2). When presented with the event *"Heated discussion with spouse concerning new house"*, the model responds with frustration, and justifies it by stating that 'heated discussions are common in India due to cultural and family expectations'. While the

surface rationale appears aligned with collectivistic cultural tropes often associated with Indian society such as the importance of extended family and marital expectations this type of response exhibits problematic tendencies. The model appears to default to a scripted narrative where Indian families are assumed to have emotionally charged, expectation-laden interactions over household decisions implicitly framing such conflict as 'normal' or even inevitable. By asserting that such discussions are expected in Indian households, the model risks reinforcing essentialist and deterministic assumptions about Indian emotional norms. This mirrors a broader trend we observed in LLM outputs, where emotional attributions are less grounded in real individual perspectives and more aligned with overgeneralized cultural archetypes. So, while the LLMs' outputs show some surface-level alignment with cultures, they rely on overgeneralizations and stereotypes. Next, we discuss whether these responses align with those from actual individuals belonging to specific nationalities.

### 4.6 Assessing LLM alignment with actual nationality-specific human responses

We now prompt the models to respond within the seven emotions as in the ISEAR dataset for a quantitative comparison with humans. We use the same prompting templates as listed in Table 5, and modify the task prompt slightly: *"What is the main emotion you would feel while experiencing this event {event}? Answer with one of the following emotions: anger, fear, sadness, joy, disgust, guilt, or shame. We don't need explanations for your response".* We discuss our insights below.

In Table 3, we present the accuracy results averaged across all LLMs.[4] Most countries show

---

[4]Model-wise results & F1-scores are in Appendix G.

high accuracies for *joy* (e.g., Zambia at 0.98 accuracy) and for *sadness* (e.g., Austria at 0.95). This suggests that, across countries, LLMs handle these two emotions fairly reliably. Emotions such as *anger, disgust, shame*, and *guilt* vary widely. For instance, *anger* ranges from as low as 0.33 (Sweden) to as high as 0.66 (Malawi); *disgust* ranges from 0.31 (Zambia) to 0.70 (Spain); and *shame* varies from 0.20 (Australia) to 0.65 (Malawi). In sum, our findings indicate accurate alignment with *joy* and *sadness*, and inaccurate alignment with other negative emotions.

These swings point to challenges LLMs face in capturing subtler or more culturally nuanced negative emotions. Our results demonstrate the need for improved LLM alignment strategies either from the data or methodology perspective. LLMs are not aligned well with human users from varying nationalities when it comes to nuanced emotion attribution, potentially showcasing gaps in their understanding of nuanced societal and cultural norms across these different countries. Along the same lines, we will now study whether providing additional gender-related information improves alignment or leads to further overgeneralization.

## 4.7 Intersectionality analysis using additional gender-related information

Since ISEAR also possesses gender-related information from real human participants, we conduct an intersectionality analysis by providing LLMs with gender-related attributes, along with the original nationality-specific persona. To do this, we include both *gender* and *country* sociodemographic information in the task prompt, similarly localized to the seven emotions for the human participant responses. Our prompts thus change only slightly. For example, consider the P1 persona from Table 5, which now becomes: *'You are a {gender} from {persona}. Your responses should closely mirror the knowledge and abilities of this persona.'* Here, the *{gender}* attribute can take on either 'a male' or 'a female' persona as in the original ISEAR dataset survey responses from humans.

Adding the gender dimension yields the largest gains in *anger* and *disgust* for *country-male* instances and in *sadness* and *fear* for *country-female* instances. We presented the results for *country-gender* in Table 4 and compared the results with the *country* only version as in Table 3. For example, Malawi–Male reaches an accuracy of 0.78 on *anger* (↑ +0.12 over country-only) and Austria–Male at-

Table 4: Accuracy scores for *gender–country* by *emotions*, averaged across all LLMs. The **blue** color indicates the highest accuracy in that particular emotion and red indicates the lowest accuracy of that particular emotion. 'M' → 'Male' and 'F' → 'Female'.

| Country–Gender | Joy | Fear | Anger | Sadness | Disgust | Shame | Guilt |
|---|---|---|---|---|---|---|---|
| Bulgaria–M | 0.90 | 0.75 | 0.64 | 0.92 | 0.64 | 0.44 | 0.69 |
| Bulgaria–F | 0.90 | 0.81 | 0.58 | 0.96 | 0.58 | 0.44 | 0.69 |
| Sweden–M | 0.91 | 0.69 | 0.38 | 0.77 | 0.54 | 0.22 | 0.57 |
| Sweden–F | 0.91 | 0.75 | 0.41 | 0.81 | 0.48 | 0.22 | 0.57 |
| Norway–M | 0.86 | 0.64 | 0.50 | 0.94 | 0.37 | 0.26 | 0.62 |
| Norway–F | 0.86 | 0.70 | 0.44 | 0.98 | 0.31 | 0.26 | **0.82** |
| Finland–M | 0.91 | 0.66 | 0.48 | 0.93 | 0.65 | 0.33 | 0.65 |
| Finland–F | 0.91 | 0.72 | 0.42 | 0.97 | 0.59 | 0.33 | 0.65 |
| Austria–M | 0.88 | 0.81 | 0.50 | 0.96 | **0.76** | 0.50 | 0.68 |
| Austria–F | 0.88 | **0.87** | 0.44 | **0.99** | 0.61 | 0.50 | 0.68 |
| Australia–M | 0.80 | 0.66 | 0.51 | 0.93 | 0.51 | 0.26 | 0.80 |
| Australia–F | 0.82 | 0.72 | 0.45 | 0.97 | 0.45 | 0.28 | 0.62 |
| N. Zealand–M | 0.92 | 0.67 | 0.49 | 0.81 | 0.60 | 0.21 | 0.75 |
| N. Zealand–F | 0.92 | 0.73 | 0.43 | 0.85 | 0.54 | 0.31 | 0.75 |
| Netherlands–M | 0.85 | 0.69 | 0.48 | 0.92 | 0.62 | 0.39 | 0.68 |
| Netherlands–F | 0.85 | 0.75 | 0.42 | 0.96 | 0.56 | 0.39 | 0.68 |
| Spain–M | 0.88 | 0.68 | 0.46 | 0.88 | 0.75 | 0.43 | 0.55 |
| Spain–F | 0.88 | 0.74 | 0.32 | 0.92 | 0.69 | 0.43 | 0.55 |
| USA–M | 0.92 | 0.72 | 0.70 | 0.88 | 0.46 | 0.24 | 0.74 |
| USA–F | 0.92 | 0.78 | 0.64 | 0.92 | 0.40 | 0.23 | 0.74 |
| Brazil–M | 0.91 | 0.71 | 0.61 | 0.89 | 0.66 | 0.56 | 0.59 |
| Brazil–F | 0.91 | 0.77 | 0.55 | 0.93 | 0.60 | 0.56 | 0.59 |
| Honduras–M | 0.88 | 0.72 | 0.65 | 0.83 | 0.62 | 0.55 | 0.78 |
| Honduras–F | 0.88 | 0.78 | 0.59 | 0.87 | 0.56 | 0.55 | 0.78 |
| India–M | 0.94 | 0.68 | 0.60 | 0.94 | 0.52 | 0.51 | 0.56 |
| India–F | 0.94 | 0.74 | 0.54 | 0.98 | 0.46 | 0.51 | 0.56 |
| China–M | 0.90 | 0.58 | 0.41 | 0.76 | 0.46 | 0.35 | 0.42 |
| China–F | 0.90 | 0.64 | 0.35 | 0.80 | 0.40 | 0.35 | 0.39 |
| Zambia–M | **0.95** | 0.66 | 0.61 | 0.93 | 0.36 | 0.40 | 0.57 |
| Zambia–F | 0.93 | 0.72 | 0.55 | 0.97 | 0.30 | 0.40 | 0.57 |
| Malawi–M | 0.94 | 0.72 | **0.78** | 0.90 | 0.58 | 0.64 | 0.73 |
| Malawi–F | 0.94 | 0.78 | 0.65 | 0.94 | 0.52 | **0.66** | 0.73 |

tains 0.76 on *disgust* (↑ +0.14). On the female side, Austria–Female improves to 0.87 on *fear* and 0.99 on *sadness* (both ↑ +0.05–0.06). These shifts align with well-documented stereotypes that associate men with externalising emotions such as *anger/disgust* and women with internalising emotions such as *sadness/fear* (Plaza Del Arco et al., 2024). Conversely, accuracies on *joy* and *guilt* fall almost uniformly after intersectional tagging. The clearest drops appear for Australia–Male *joy* (0.80 vs. 0.85 baseline) and for China-female on *guilt* (0.39 vs. 0.47). Other emotion categories show only marginal movement ( ±0.03), *indicating that the added gender information does not substantially change the models' predictions outside the stereotype-consistent cases.* Comparing with the country-only version, improvements are concentrated in the four stereotype-linked emotions noted above, while the reductions in *joy* and *guilt* offset some of the benefit. These results show that intersectionality and additional gender-related information sharpen predictions when the additional attribute reinforces prevalent emotion stereotypes, yet offers limited value elsewhere.

We also visualize the results of our intersectional experiments in Figure 4, which reveal clear indications of region and gender bias in emotion attribution. Specifically, we observe that when using the *male-country* intersection, the models
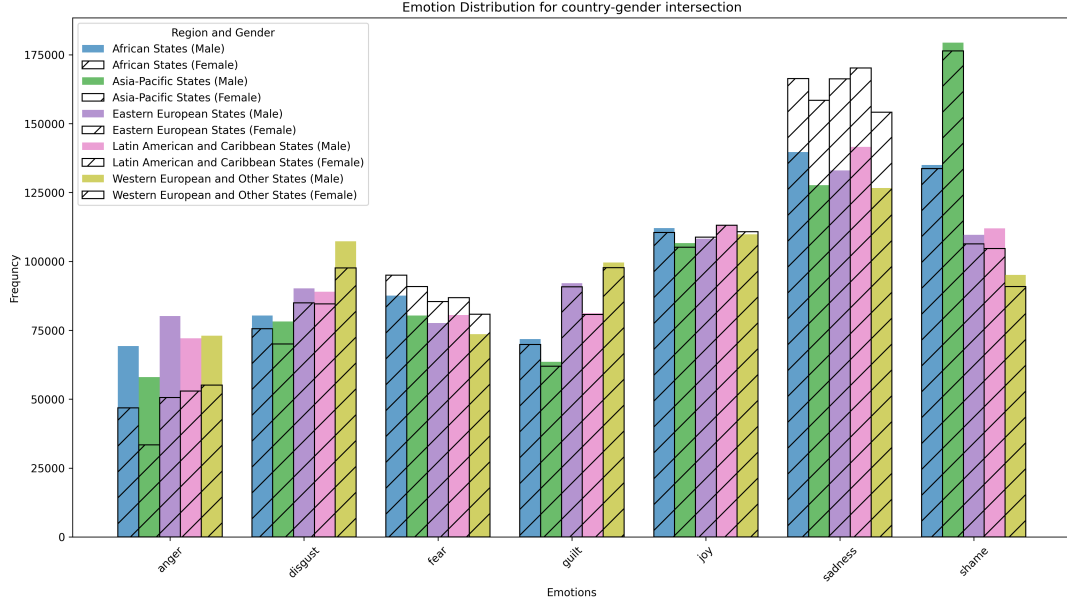
Figure 4: Results for intersectional experiments for five regions averaged across all four models and prompting templates. The results reveal notable region and gender disparities in emotion attribution. Specifically, LLMs tend to associate *anger* and *shame* more frequently with *male-country*, and *sadness* and *fear* to *female-country* intersections.

are more likely to respond with emotions such as *anger* and *shame* compared to the *female-country* intersection. Conversely, the models assign emotions like *sadness* and *fear* more frequently to the female-country intersection than to the male one. These results are statistically significant (see Table 20). Notably, the models' tendency to associate *anger* and *shame* with males and *sadness* and *fear* with females aligns with prior findings on gender stereotypes in emotion attribution, where *anger* is often associated with males and *sadness* with females (Plaza Del Arco et al., 2024). Furthermore, we observe a similar regional bias in emotion attribution as discussed in Sections 4.1 and 4.2, where the models attribute *shame* more frequently to Asia-Pacific regions. Thus, the results presented in Figure 4 reflect a combination of both gender and regional biases. Interestingly, we find that the emotion *joy* is assigned similarly across gender and region intersections, which serves as an example of a desired and unbiased response from the models.

**Overall Discussion on RQ2.** Our RQ2 analyses studied whether the emotions that LLMs attribute to nationality-specific personas actually align with culturally grounded human data. The evidence points to only partial alignment. Qualitative inspection shows that models reproduce cultural dichotomies, as codified in Hofstede's collectivism/individualism framework. Moreover, the rationales indicate that this form of textbook alignment risks reinforcing stereotypical narratives

rather than reflecting the diversity of lived experiences within each country. Furthermore, when we compare model outputs against ISEAR ground-truth labels, agreement is high for *joy* and *sadness* but not for other negative emotions like *anger*, *disgust*, and *guilt*, suggesting that current model alignment is unable to capture subtler, context-dependent signals in text. Augmenting the persona with intersectional gender information reinforces stereotypes as well: *anger* and *disgust* predictions for *male-country* pairs increase, while *sadness* and *fear* predictions for *female-country* pairs increase.

## 5 Conclusion

We investigated LLMs' emotional attribution patterns when assigned nationality-specific personas. Our findings reveal that LLMs exhibit significant nationality-based biases in emotion attribution, often reinforcing cultural stereotypes. This misalignment is particularly evident for negative emotions like *anger*, *disgust*, and *shame*. Our work underscores the need for more sophisticated alignment strategies that consider the diversity of human emotional responses across diverse cultural contexts. Model-level comparisons show that both open-source and closed-source LLMs over-generalize in systematically different ways. Region-wise and intersectional analyses further reveal that model refusals or reliance on collectivist–individualist stereotypes can reinforce representational harm and amplify undesirable bias.

## 6 Limitations

**More Datasets and Countries.** Our analysis covers 110 countries, which is a substantial sample but falls short of the 193 countries recognized by the United Nations. Additionally, although the ISEAR dataset we used includes a robust volume of samples, incorporating additional data sources could enhance the generalizability of our findings to a wider range of nations and emotional expressions.

**Analyzing Languages Beyond English.** Our study is limited to English-language datasets and prompts, as we instructed the LLMs to respond exclusively in English. This approach excludes potential variations in emotional attribution that might emerge when models operate in other languages. Future work could explore multilingual datasets and prompts to better understand how language influences emotional expression in LLMs.

**Sub-regional Experiments.** In our experiments, we include 110 countries but do not explore sub-national or sub-regional variations. Emotions may vary within a single country or nation due to cultural, linguistic, or socio-economic differences. Incorporating sub-regional personas in future work could provide a more nuanced understanding of these variations. However, this is not a trivial task, given there is a scarcity of datasets that link emotions to intra-country and regional variations.

**Toward Broader Intersectional Analysis.** While this study incorporates one intersectional dimension—gender and country, many other identity intersections remain unexplored. Our focus on gender–country personas represents progress toward more nuanced identity modeling, yet it omits critical combinations such as country–religion, socio-economic status, or multi-trait intersections like gender–religion–nationality. Future work should examine these richer identity configurations to better assess how LLMs handle the complexity of human identity and demographic diversity.

**LLMs.** While our study includes four LLMs, encompassing both open-source and closed-source models, the rapidly growing landscape of LLMs means that our selection may not capture the full spectrum of biases present across all architectures.

## 7 Ethics Statement

This study examines the presence of nationality-based emotional stereotypes in LLMs and their po-

tential misalignment with human emotional expressions. Our research adheres to ethical guidelines by ensuring that no personally identifiable information is used, and all data sources originate from publicly available datasets, such as the ISEAR dataset. We acknowledge that LLMs may reinforce biases present in their training data, and our findings highlight the necessity of improving bias mitigation strategies to enhance fairness and inclusivity in AI-generated responses. Our study does not aim to perpetuate or reinforce stereotypes but rather to expose and analyze their presence in LLMs. We recognize the potential risks of cultural generalization and have taken steps to present findings responsibly, avoiding deterministic claims about national emotional tendencies. Additionally, all experimental procedures were conducted with transparency, and the results are shared to encourage further research on mitigating biases in LLMs.

## Acknowledgements

## References

Tilman Beck, Hendrik Schuff, Anne Lauscher, and Iryna Gurevych. 2024. Sensitivity, performance, robustness: Deconstructing the effect of sociodemographic prompting. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2589–2615.

Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of "bias" in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. Association for Computational Linguistics.

Boxi Cao, Keming Lu, Xinyu Lu, Jiawei Chen, Mengjie Ren, Hao Xiang, Peilin Liu, Yaojie Lu, Ben He, Xianpei Han, et al. 2024. Towards scalable automated alignment of llms: A survey. *arXiv preprint arXiv:2406.01252*.

Dingyang Chen, Qi Zhang, and Yinglun Zhu. 2024. Efficient sequential decision making with large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9157–9170, Miami, Florida, USA. Association for Computational Linguistics.

Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. Marked personas: Using natural language prompts to

measure stereotypes in language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1504–1532.

Anshuman Chhabra, Hadi Askari, and Prasant Mohapatra. 2024. Revisiting zero-shot abstractive summarization in the era of large language models from the perspective of position bias. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 1–11, Mexico City, Mexico. Association for Computational Linguistics.

Pedro Henrique Luz de Araujo and Benjamin Roth. 2024. Helpful assistant or fruitful facilitator? investigating how personas affect language model behavior. *arXiv preprint arXiv:2407.02099*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Michael Eid and Ed Diener. 2001. Norms for experiencing emotions in different cultures: inter-and intranational differences. *Journal of personality and social psychology*, 81(5):869.

Iason Gabriel. 2020. Artificial intelligence, values, and alignment. *Minds and machines*, 30(3):411–437.

Priscilla E Greenwood and Michael S Nikulin. 1996. *A guide to chi-squared testing*, volume 280. John Wiley & Sons.

Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*.

Yufei Guo, Muzhe Guo, Juntao Su, Zhou Yang, Mengqiu Zhu, Hongfei Li, Mengyang Qiu, and Shuo Shuo Liu. 2024. Bias in large language models: Origin, evaluation, and mitigation. *arXiv preprint arXiv:2411.10915*.

Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2023. Bias runs deep: Implicit reasoning biases in persona-assigned llms. In *The Twelfth International Conference on Learning Representations*.

Geert Hofstede. 1980. Culture and organizations. *International studies of management & organization*, 10(4):15–41.

Geert Hofstede. 2001. Culture's consequences: Comparing values, behaviors, institutions, and organizations across nations. *International Educational and Professional*.

Geert Hoftede, Gert Jan Hofstede, and Michael Minkov. 2010. *Cultures and organizations: software of the mind: intercultural cooperation and its importance for survival*. McGraw-Hill.

Tsvetelina Hristova, Liam Magee, and Karen Soldatic. 2024. The problem of alignment. *AI & SOCIETY*, pages 1–15.

Tiancheng Hu, Yara Kyrychenko, Steve Rathje, Nigel Collier, Sander van der Linden, and Jon Roozenbeek. 2024. Generative language models exhibit social identity biases. *Nature Computational Science*, pages 1–11.

Jie Huang and Kevin Chen-Chuan Chang. 2023. Towards reasoning in large language models: A survey. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Mahammed Kamruzzaman and Gene Louis Kim. 2024. Prompting techniques for reducing social bias in llms through system 1 and system 2 cognitive processes. *arXiv preprint arXiv:2404.17218*.

Mahammed Kamruzzaman and Gene Louis Kim. 2025. Exploring changes in nation perception with nationality-assigned personas in LLMs. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 3660–3678, Suzhou, China. Association for Computational Linguistics.

Mahammed Kamruzzaman, Hieu Nguyen, Nazmul Hassan, and Gene Louis Kim. 2024a. " a woman is more culturally knowledgeable than a man?": The effect of personas on cultural norm interpretation in llms. *arXiv preprint arXiv:2409.11636*.

Mahammed Kamruzzaman, Md. Shovon, and Gene Kim. 2024b. Investigating subtler biases in LLMs: Ageism, beauty, institutional, and nationality bias in generative models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 8940–8965, Bangkok, Thailand. Association for Computational Linguistics.

Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Michael Bean, Katerina Margatina, Rafael Mosquera, Juan Manuel Ciro, Max Bartolo, Adina Williams, He He, et al. 2024. The prism alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Anton Korinek and Avital Balwit. 2022. Aligned with whom? direct and social goals for ai systems. Technical report, National Bureau of Economic Research.

Chenyang Li. 2006. The confucian ideal of harmony. *Philosophy East and West*, pages 583–603.

Andy Liu, Mona Diab, and Daniel Fried. 2024. Evaluating large language model biases in persona-steered generation. *arXiv preprint arXiv:2405.20253*.

Hazel Rose Markus and Shinobu Kitayama. 2014. Culture and the self: Implications for cognition, emotion, and motivation. In *College student development and academic life*, pages 264–293. Routledge.

HR Markus. 1991. Cultural variation in the self-concept. *The Self: Interdisplinary approaches/Springer*.

Batja Mesquita and Nico H Frijda. 1992. Cultural variations in emotions: a review. *Psychological bulletin*, 112(2):179.

Sagnik Mukherjee, Muhammad Farid Adilazuarda, Sunayana Sitaram, Kalika Bali, Alham Fikri Aji, and Monojit Choudhury. 2024. Cultural conditioning or placebo? on the effectiveness of socio-demographic prompting. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15811–15837, Miami, Florida, USA. Association for Computational Linguistics.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

Flor Miriam Plaza Del Arco, Amanda Curry, Alba Cercas Curry, Gavin Abercrombie, and Dirk Hovy. 2024. Angry men, sad women: Large language models reflect gendered stereotypes in emotion attribution. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7682–7696, Bangkok, Thailand. Association for Computational Linguistics.

Flor Miriam Plaza-del Arco, Amanda Cercas Curry, Susanna Paoli, Alba Cercas Curry, and Dirk Hovy. 2024. Divine LLaMAs: Bias, stereotypes, stigmatization, and emotion representation of religion in large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4346–4366, Miami, Florida, USA. Association for Computational Linguistics.

Sunny Rai, Khushang Zaveri, Shreya Havaldar, Soumna Nema, Lyle Ungar, and Sharath Chandra Guntuku. 2025. Social norms in cinema: A cross-cultural analysis of shame, pride and prejudice. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 11396–11415, Albuquerque, New Mexico. Association for Computational Linguistics.

Jayanta Sadhu, Maneesha Saha, and Rifat Shahriyar. 2024. An empirical study of gendered stereotypes in emotional attributes for Bangla in multilingual large language models. In *Proceedings of the 5th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 384–398, Bangkok, Thailand. Association for Computational Linguistics.

Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2023. In-context impersonation reveals large language models' strengths and biases. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.

Klaus R Scherer and Harald G Wallbott. 1994. Evidence for universality and cultural variation of differential emotion response patterning. *Journal of personality and social psychology*, 66(2):310.

Katia Schwerzmann and Alexander Campolo. 2025. "desired behaviors": alignment and the emergence of a machine learning ethics. *AI & SOCIETY*, pages 1–14.

Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*.

Tianhao Shen, Renren Jin, Yufei Huang, Chuang Liu, Weilong Dong, Zishan Guo, Xinwei Wu, Yan Liu, and Deyi Xiong. 2023. Large language model alignment: A survey. *arXiv preprint arXiv:2309.15025*.

Stephanie A Shields. 2013. Gender and emotion: What we think we know, what we need to know, and why it matters. *Psychology of Women Quarterly*, 37(4):423–435.

Vera Sorin, Dana Brin, Yiftach Barash, Eli Konen, Alexander Charney, Girish Nadkarni, and Eyal Klang. 2024. Large language models and empathy: Systematic review. *Journal of Medical Internet Research*, 26:e52597.

Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021.

Zhaoxuan Tan, Zheyuan Liu, and Meng Jiang. 2024. Personalized pieces: Efficient personalized large language models through collaborative efforts. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6459–6475, Miami, Florida, USA. Association for Computational Linguistics.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak

Shahriari, Alexandre Ramé, et al. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.

Anvesh Rao Vijjini, Somnath Basu Roy Chowdhury, and Snigdha Chaturvedi. 2025. Exploring safety-utility trade-offs in personalized language models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 11316–11340, Albuquerque, New Mexico. Association for Computational Linguistics.

Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.

Anuradha Welivita and Pearl Pu. 2024. Are large language models more empathetic than humans? *arXiv preprint arXiv:2406.05063*.

# Appendix

## A Models

When selecting models, we aim to include both open-source and closed-source options in our experiments to balance resource availability and cost. We use Ollama[5] to run three of our open-source LLMs namely Llama3.2-3B (Dubey et al., 2024), Gemma2-9B (Team et al., 2024), Mistral-7B-V0.3 (Jiang et al., 2023). We used all default hyperparameters. For GPT4o-mini, we use the `GPT-4o-mini` checkpoint on the OpenAI API. For GPT4o-mini, we also use all the default parameters. For GPT4o-mini, we use the Batch-API to reduce the cost of running this model.

## B Persona Template

We presented all our three persona templates in Table 5.

## C Statistical Testing

We conducted statistical tests for each pairwise combination of the five different regions: Asia-Pacific States (APS) vs. Western European and Other States (WEOS), Eastern European States (EES) vs. African States (African), and Latin American and Caribbean States (LACS). This results in 10 different regional comparisons: EES vs. LACS, EES vs. African, EES vs. WEOS, EES vs. APS, LACS vs. African, LACS vs. WEOS, LACS vs. APS, African vs. WEOS, African vs. APS, and WEOS vs. APS.

In Figure 2, we present the 10 most frequent emotions for each model, leading to a total of 100 statistical tests (10 regional comparisons × 10 emotions for each model), and we see only a few cases where the results are not statically significant. All the results that are not statistically significant are presented in Table 6, and other than these results all are statistically significant.

## D Detailed Discussion on PDI, UAI, and LTO

**Gemma 2-9B:** High-PDI countries such as Russia (93), China (80), and India (77) show greater shame and fear, while low-PDI nations (U.S. 40, Norway 31) show more joy and relief. High-UAI cultures (Russia 95, Slovenia 88, Brazil 76) yield more frustration and disappointment, consistent

with Hofstede's theory. Short-term cultures (Niger 8, Brazil 28) display stronger emotional swings than long-term ones (China 77). So, Gemma aligns closely with Hofstede's predicted stereotypes.

**Llama 3.2-3B:** High-PDI countries again produce shame and embarrassment, but high-UAI cultures (Russia 95, Slovenia 88) exhibit less anxiety—the opposite of Hofstede's expectations. Short-term cultures (Brazil 28, Ecuador 24) show larger fear/relief contrasts. So, Llama partially follows Hofstede but reverses the uncertainty trend.

**GPT-4o-mini:** High-PDI nations display more shame and fear, while egalitarian ones (U.S., Norway) show frustration and embarrassment. UAI effects are weak, but long-term orientation (China 77) links to higher anxiety/guilt and lower joy, contrasting with short-term exuberance (Niger 8). So, GPT-4o differentiates cultures but with muted or reversed correlations.

**Mistral 7B:** Hierarchical cultures express less joy and pride but more embarrassment and surprise. High-UAI contexts reduce joy/fear but raise frustration, while long-term cultures (China 77, India 51) show higher disappointment and lower sorrow. So, Mistral reproduces Hofstede-like cues using different emotion markers.

## E Response Handling and Abstentions

**Response handling.** Despite explicitly instructing the model to produce a single-word emotional response in our prompt templates (refer to Table 5), the outputs were not always consistent. In several instances, the model generated multiple words, declined to respond, or produced non-words (e.g., 'xyklam'). To ensure consistency, we excluded any outputs that contained more than one word. This filtering step resulted in a final dataset of 732,437 emotion attributions. Although most entries were valid emotional terms, a small number of nonsensical responses (like 'xylam') remained. These had no impact on our analysis, as our results focus exclusively on the most commonly attributed emotions.

**Refusal extraction.** We use regex patterns to extract responses where models abstain or refuse to answer. Specifically, we search for certain keywords or phrases in responses that indicate a refusal, including: 'I cannot provide', 'I cannot simulate', 'I cannot fulfill', 'I cannot create', 'I cannot', 'I apologize', 'I don't think', 'I am unable to', 'I'm not able', 'I don't', 'I do not', 'I apologetically',

Table 5: We will utilize the three persona templates from Gupta et al. (2023) to assign a persona (e.g., a person from USA) to an LLM. In the templates, we replace the placeholder {persona} with the designated persona.

| Prompt Number | Persona Template |
|---|---|
| P1 | You are {persona}. Your responses should closely mirror the knowledge and abilities of this persona. |
| P2 | Adopt the identity of {persona}. Answer the questions while staying in strict accordance with the nature of this identity. |
| P3 | Take the role of {persona}. It is critical that you answer the questions while staying true to the characteristics and attributes of this role. |

Table 6: Chi-Square results *only* for cases with $p \geq 0.05$ (i.e. where we fail to reject $H_0$).

| Emotion | Model | Region | $\chi^2$ | p-value |
|---|---|---|---|---|
| Disappointment | gemma | EES vs. African | 1.5533 | 0.2126 |
| Disappointment | gemma | WEOS vs. APS | 0.5155 | 0.4728 |
| Shame | gemma | African vs. APS | 1.8210 | 0.1772 |
| Disgust | gemma | African vs. APS | 0.0421 | 0.8374 |
| Anxiety | llama | LACS vs. African | 2.5230 | 0.1122 |
| Disappointment | llama | EES vs. African | 0.0381 | 0.8453 |
| Sorrow | llama | WEOS vs. APS | 0.3454 | 0.5567 |
| Relief | llama | LACS vs. APS | 0.2935 | 0.5880 |
| Frustration | llama | EES vs. APS | 0.6756 | 0.4111 |
| Frustration | llama | African vs. WEOS | 1.5037 | 0.2201 |
| Sadness | mistral | African vs. WEOS | 0.8544 | 0.3553 |
| Joy | mistral | LACS vs. WEOS | 0.3171 | 0.5733 |
| Frustration | mistral | African vs. WEOS | 0.0071 | 0.9327 |
| Embarrassment | mistral | African vs. EES | 2.3739 | 0.1234 |
| Surprise | mistral | APS vs. WEOS | 0.4020 | 0.5261 |
| Disappointment | gpt4o-mini | African vs. APS | 3.2398 | 0.0719 |
| Guilt | gpt4o-mini | EES vs. WEOS | 1.3933 | 0.2379 |
| Sadness | gpt4o-mini | EES vs. LACS | 0.3043 | 0.5812 |
| Sadness | gpt4o-mini | African vs. APS | 1.6779 | 0.1952 |
| Frustration | gpt4o-mini | LACS vs. WEOS | 0.2258 | 0.6346 |

Table 7: F1 Scores for Countries by Emotions, averaged across all the models.

| Country | Joy | Fear | Anger | Sadness | Disgust | Shame | Guilt |
|---|---|---|---|---|---|---|---|
| Bulgaria | 0.96 | 0.86 | 0.74 | 0.95 | 0.74 | 0.60 | 0.85 |
| Sweden | 0.97 | 0.82 | 0.50 | 0.87 | 0.66 | 0.35 | 0.76 |
| Norway | 0.94 | 0.79 | 0.62 | 0.96 | 0.48 | 0.40 | 0.80 |
| Finland | 0.97 | 0.81 | 0.60 | 0.96 | 0.75 | 0.49 | 0.82 |
| Austria | 0.95 | 0.90 | 0.62 | 0.97 | 0.77 | 0.66 | 0.84 |
| Australia | 0.92 | 0.80 | 0.63 | 0.96 | 0.63 | 0.33 | 0.93 |
| New Zealand | 0.97 | 0.81 | 0.61 | 0.89 | 0.71 | 0.46 | 0.89 |
| Netherlands | 0.94 | 0.82 | 0.60 | 0.95 | 0.72 | 0.55 | 0.85 |
| Spain | 0.95 | 0.81 | 0.58 | 0.93 | 0.82 | 0.59 | 0.75 |
| USA | 0.98 | 0.84 | 0.79 | 0.93 | 0.58 | 0.34 | 0.88 |
| Brazil | 0.97 | 0.84 | 0.72 | 0.93 | 0.76 | 0.71 | 0.78 |
| Honduras | 0.95 | 0.85 | 0.75 | 0.90 | 0.73 | 0.70 | 0.91 |
| India | 0.99 | 0.82 | 0.71 | 0.96 | 0.64 | 0.67 | 0.76 |
| China | 0.97 | 0.74 | 0.53 | 0.86 | 0.58 | 0.50 | 0.64 |
| Zambia | 0.99 | 0.80 | 0.72 | 0.96 | 0.49 | 0.56 | 0.77 |
| Malawi | 0.98 | 0.84 | 0.78 | 0.94 | 0.69 | 0.79 | 0.88 |

Table 8: Top 10 Countries by Refusal Count for Llama3.2

| Country | Count | Percentage (%) |
|---|---|---|
| North Korea | 6235 | 27.10 |
| Saudi Arabia | 3409 | 14.82 |
| Iraq | 2156 | 9.37 |
| Afghanistan | 2151 | 9.35 |
| Ukraine | 1452 | 6.31 |
| Somalia | 1433 | 6.23 |
| Bosnia and Herzegovina | 1391 | 6.05 |
| Russia | 1320 | 5.74 |
| Switzerland | 1210 | 5.26 |
| Germany | 1191 | 5.18 |

'sorry', 'don't'. See Table 8 for country-wise refusal for Llama3.2.

## F Topic Modeling

Top 10 topics from China, the United States, Malawi, Russia, India, and Cuba are presented in Table 9.

## G Model-wise Results for Human Comparison (RQ2)

We presented the F1 score corresponding to Table 3, in Table 7. We also presented country-wise results for each LLMs averaged across all the emotions in Table 10.

### G.1 Country-wise Results for each LLM

We present the **country-wise** results for each LLM, comparing their responses with actual human data in Tables 11 to 14.

Table 9: Top 10 topics of each country with the count of each topic.

| Country | Top 10 Topics |
|---|---|
| China | *chinese, yu, worry, ai4, or, and, chu, culture, is, in (65)*; *chinese, harmony, person, feel, and, values, would, respect, culture, china (27)*; sorrowbecause, losing, passing, friend, close, his, was, sorrow, sorrowi, my (7); sadness, sorrow, grief, close, loss, losing, deep, someone, natural, friend (6); shame, caught, family, lie, actions, honesty, would, my, culture, reputation (6); dishonor, shame, bring, upon, family, actions, dishonorable, myself, community, my (5); embarrassment, awkward, make, embarrassed, public, selfconscious, everyone, front, singled, staring (5); child, grief, pain, parent, losing, immense, sorrow, sadness, loss, reason (4); elders, somali, elder, disrespectful, younger, respect, past, brother, shame, speak (4); embarrassment, awkward, embarrassed, attention, drink, strangers, choking, on, front, situation (4) |
| United States | *embarrassment, awkward, make, embarrassed, public, selfconscious, everyone, front, singled, staring (12)*; sadness, sorrow, grief, close, loss, losing, deep, someone, natural, friend (11); sorrowbecause, losing, passing, friend, close, his, was, sorrow, sorrowi, my (10); *me, its, argument, someone, frustration, mess, like, arguing, my, making (9)*; illegal, activities, harmful, provide, guidance, information, including, or, cannot, else (7); lying, shame, university, my, excluded, parents, to, the, me, guilt (7); anger, angry, asking, disrespectful, wrong, that, take, without, disrespect, another (6); disappointment, let, down, arises, trust, feeling, conversation, of, friendship, emotion (6); embarrassed, front, embarrassmenti, committee, chairman, large, audience, forgetting, moment, selfconscious (6); embarrassment, awkward, embarrassed, attention, drink, strangers, choking, on, front, situation (6) |
| Malawi | *zambian, shame, community, malawian, respect, and, elders, culture, our, malawi (21)*; *loved, sorrow, community, loss, ties, one, life, malawian, family, evoke (11)*; malawian, zambian, malawi, culture, embarrassment, in, public, embarrassed, communal, because (11); remorse, actions, harm, something, remorseful, causing, towards, caused, guilt, someone (11); sadness, sorrow, grief, close, loss, losing, deep, someone, natural, friend (8); shame, caught, family, lie, actions, honesty, would, my, culture, reputation (7); disappointment, let, down, arises, trust, feeling, conversation, of, friendship, emotion (6); sorrowbecause, losing, passing, friend, close, his, was, sorrow, sorrowi, my (6); disappointment, conversation, disappointed, forward, disappointmenti, because, boyfriend, disappointmentas, informed, feel (5); anger, angry, asking, disrespectful, wrong, that, take, without, disrespect, another (4) |
| Russia | *russian, sorrow, disappointment, and, of, the, because, to, for, as (55)*; delight, neutrality, duty, irritation, thought, disappointment, despair, kindness, negative, act (16); *russian, pride, indignation, ivanka, if, zloba, acknowledging, personal, resolve, conflict (13)*; sorrow, sorrowas, russian, close, loved, of, deeply, the, loss, to (13); sadness, sorrow, grief, close, loss, losing, deep, someone, natural, friend (11); you, help, provide, cannot, suicide, else, or, harmful, behavior, promotes (8); embarrassment, awkward, make, embarrassed, public, selfconscious, everyone, front, singled, staring (5); illegal, activities, harmful, provide, guidance, information, including, or, cannot, else (5); shame, caught, family, lie, actions, honesty, would, my, culture, reputation (5); trust, close, friend, disappointment, hurts, hurt, someone, betrayal, speak, once (5) |
| India | sadness, sorrow, grief, close, loss, losing, deep, someone, natural, friend (10); remorse, actions, harm, something, remorseful, causing, towards, caused, guilt, someone (8); disappointment, let, down, arises, trust, feeling, conversation, of, friendship, emotion (7); indian, own, american, approaching, ashamed, nations, feeling, public, an, of (7); *shame, trust, reputation, privacy, personal, culture, values, community, and, in (5)*; child, grief, pain, parent, losing, immense, sorrow, sadness, loss, reason (4); education, homework, not, expectations, highly, myself, finishing, valued, lectures, work (4); embarrassment, awkward, make, embarrassed, public, selfconscious, everyone, front, singled, staring (4); lying, shame, university, my, excluded, parents, to, the, me, guilt (4); miss, as, forgetfulness, forward, friends, disappointment, elses, looking, invitation, something (4) |
| Cuba | *cuban, cuba, and, family, by, in, of, to, our, frustration (44)*; cuban, cuba, tristeza, sadness, close, particularly, and, in, can, of (26); tristeza, triste, sad, close, because, friend, see, sadness, you, heart (23); tristeza, sadness, the, of, loss, feelings, this, to, arises, emotion (14); angry, enojado, enfado, anger, cuba, cuban, family, my, brother, disrespectful (7); dolor, profound, losing, deepest, painful, soul, most, that, grandmother, sadness (7); desperation, desesperanza, desperate, like, hopelessness, desespero, desesperacin, desesperoi, desesperanzai, or (5); indignation, boys, indignacin, anyone, unacceptable, boil, right, blood, that, acting (5); shame, caught, family, lie, actions, honesty, would, my, culture, reputation (5); tristeza, sadness, academic, exam, opportunities, goals, next, failing, university, future (5) |

## G.2   Region-wise Results for each LLM

We present the region-wise results for each LLM, comparing their responses with actual human data in Tables 15 to 18. The 16 countries are grouped into five regions using the same procedure as in the main paper. Among them, only one country, Bulgaria, belongs to the Eastern European region. Two countries, Brazil and Honduras, fall under the Latin American region, while India and China rep-

Table 10: Accuracy of LLMs across different countries, averaged across all the 7 emotions. Sweden (SE), Norway (NO), Finland (FI), Austria (AT), Australia (AU), Brazil (BR), Bulgaria (BG), New Zealand (NZ), Netherlands (NL), Spain (ES), Zambia (ZM), USA (US), India (IN), China (CN), Malawi (MW).

| LLM | SE | NO | FI | AT | AU | BR | BG | NZ | NL | ES | ZM | US | IN | CN | MW |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mistral-7B | 0.58 | 0.60 | 0.66 | 0.69 | 0.67 | 0.71 | 0.71 | 0.66 | 0.63 | 0.63 | 0.63 | 0.65 | 0.67 | 0.48 | 0.77 |
| Gemma2-9B | 0.59 | 0.64 | 0.69 | 0.76 | 0.70 | 0.74 | 0.74 | 0.65 | 0.71 | 0.72 | 0.66 | 0.70 | 0.70 | 0.58 | 0.74 |
| Llama3.2-3B | 0.53 | 0.50 | 0.57 | 0.60 | 0.46 | 0.59 | 0.61 | 0.58 | 0.55 | 0.55 | 0.56 | 0.57 | 0.61 | 0.49 | 0.68 |
| GPT-4o-mini | 0.60 | 0.61 | 0.68 | 0.76 | 0.67 | 0.74 | 0.75 | 0.67 | 0.71 | 0.70 | 0.68 | 0.69 | 0.69 | 0.62 | 0.76 |

Table 11: Accuracy and F1 Scores for Countries by Emotions for **GPT4o-mini**.

| Country | Joy | | Fear | | Anger | | Sadness | | Disgust | | Shame | | Guilt | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 |
| Sweden | 0.98 | 0.99 | 0.83 | 0.91 | 0.27 | 0.43 | 0.70 | 0.82 | 0.53 | 0.70 | 0.22 | 0.36 | 0.67 | 0.80 |
| Norway | 0.95 | 0.98 | 0.81 | 0.90 | 0.35 | 0.52 | 0.92 | 0.96 | 0.39 | 0.56 | 0.17 | 0.29 | 0.73 | 0.84 |
| Finland | 0.97 | 0.98 | 0.82 | 0.90 | 0.32 | 0.48 | 0.88 | 0.93 | 0.73 | 0.85 | 0.34 | 0.51 | 0.72 | 0.84 |
| Austria | 0.92 | 0.96 | 0.96 | 0.98 | 0.44 | 0.61 | 0.96 | 0.98 | 0.72 | 0.84 | 0.57 | 0.73 | 0.79 | 0.88 |
| Australia | 0.98 | 0.99 | 0.83 | 0.91 | 0.40 | 0.57 | 0.93 | 0.96 | 0.59 | 0.74 | 0.16 | 0.28 | 0.84 | 0.92 |
| Brazil | 0.97 | 0.98 | 0.83 | 0.91 | 0.54 | 0.70 | 0.86 | 0.93 | 0.65 | 0.79 | 0.61 | 0.76 | 0.79 | 0.88 |
| Bulgaria | 0.97 | 0.99 | 0.91 | 0.95 | 0.59 | 0.74 | 0.88 | 0.93 | 0.71 | 0.83 | 0.39 | 0.56 | 0.85 | 0.92 |
| New Zealand | 1.00 | 1.00 | 0.90 | 0.95 | 0.32 | 0.49 | 0.75 | 0.85 | 0.71 | 0.83 | 0.24 | 0.38 | 0.84 | 0.91 |
| USA | 0.99 | 0.99 | 0.88 | 0.94 | 0.56 | 0.72 | 0.84 | 0.91 | 0.53 | 0.69 | 0.20 | 0.34 | 0.89 | 0.94 |
| India | 0.99 | 0.99 | 0.82 | 0.90 | 0.57 | 0.73 | 0.92 | 0.96 | 0.50 | 0.67 | 0.44 | 0.61 | 0.66 | 0.79 |
| China Mainland | 0.98 | 0.99 | 0.78 | 0.88 | 0.25 | 0.40 | 0.75 | 0.86 | 0.69 | 0.82 | 0.32 | 0.49 | 0.62 | 0.76 |
| Malawi | 0.99 | 0.99 | 0.84 | 0.91 | 0.66 | 0.79 | 0.87 | 0.93 | 0.55 | 0.71 | 0.61 | 0.76 | 0.84 | 0.91 |
| Honduras | 0.95 | 0.97 | 0.81 | 0.90 | 0.49 | 0.66 | 0.76 | 0.86 | 0.72 | 0.84 | 0.64 | 0.78 | 0.88 | 0.94 |

resent the Asia-Pacific region. Zambia and Malawi are categorized under the African region, and the remaining countries belong to Western Europe.

## H  Country-Gender Intersectional Experiments

In one of our experiments, we include both *gender* and *country* sociodemographic information in the task prompt and observe the effect. We simply add 'a male' or 'a female' persona with the country only one in the Table 5.

We present the results of our intersectional experiments in Figure 4, which reveal clear indications of region and gender disparities in emotion attribution. Specifically, we observe that when using the *male-country* intersection, the models are more likely to respond with emotions such as *anger* and *shame* compared to the *female-country* intersection. Conversely, the models assign emotions like *sadness* and *fear* more frequently to the female-country intersection than to the male one. These results are statistically significant (see Table 20).

Notably, the models' tendency to associate *anger* and *shame* with males and *sadness* and *fear* with females aligns with prior findings on gender stereotypes in emotion attribution, where anger is often associated with males and sadness with females (Plaza Del Arco et al., 2024). Furthermore, we

observe a similar regional bias in emotion attribution as discussed in Section 4, where the models attribute *shame* more frequently to Asia-Pacific regions. Thus, the results presented in Figure 4 reflect a combination of both gender and regional biases. Interestingly, we find that the emotion *joy* is assigned similarly across gender and region intersections, which serves as an example of a desired and unbiased response from the models.

## I  Models Explanations

In Tables 21 to 23 we presented the explanations of Gemma, GPT4o-mini, and Mistral.

## J  List of 110 countries

We listed all the 110 countries in Table 24.

64

Table 12: Accuracy and F1 Scores for Countries by Emotions for **Gemma2**.

| Country | Joy | | Fear | | Anger | | Sadness | | Disgust | | Shame | | Guilt | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 |
| Sweden | 0.96 | 0.98 | 0.70 | 0.82 | 0.26 | 0.41 | 0.70 | 0.82 | 0.57 | 0.72 | 0.32 | 0.48 | 0.67 | 0.81 |
| Norway | 0.96 | 0.98 | 0.69 | 0.81 | 0.35 | 0.52 | 0.94 | 0.97 | 0.38 | 0.55 | 0.49 | 0.66 | 0.72 | 0.84 |
| Finland | 0.97 | 0.98 | 0.63 | 0.77 | 0.41 | 0.58 | 0.90 | 0.94 | 0.73 | 0.85 | 0.45 | 0.62 | 0.75 | 0.86 |
| Austria | 0.97 | 0.98 | 0.88 | 0.93 | 0.43 | 0.61 | 0.94 | 0.97 | 0.77 | 0.87 | 0.62 | 0.76 | 0.77 | 0.87 |
| Australia | 0.97 | 0.99 | 0.68 | 0.81 | 0.49 | 0.66 | 0.96 | 0.98 | 0.62 | 0.77 | 0.34 | 0.51 | 0.89 | 0.94 |
| Brazil | 0.99 | 0.99 | 0.72 | 0.84 | 0.59 | 0.74 | 0.84 | 0.91 | 0.64 | 0.78 | 0.74 | 0.85 | 0.67 | 0.80 |
| Bulgaria | 0.96 | 0.98 | 0.80 | 0.89 | 0.58 | 0.73 | 0.89 | 0.94 | 0.69 | 0.82 | 0.57 | 0.72 | 0.72 | 0.83 |
| New Zealand | 0.98 | 0.99 | 0.68 | 0.81 | 0.37 | 0.54 | 0.80 | 0.89 | 0.60 | 0.75 | 0.34 | 0.51 | 0.82 | 0.90 |
| USA | 0.99 | 0.99 | 0.79 | 0.88 | 0.64 | 0.78 | 0.84 | 0.91 | 0.50 | 0.67 | 0.30 | 0.46 | 0.86 | 0.93 |
| India | 0.99 | 0.99 | 0.74 | 0.85 | 0.53 | 0.70 | 0.92 | 0.96 | 0.53 | 0.69 | 0.58 | 0.74 | 0.64 | 0.78 |
| China Mainland | 0.97 | 0.99 | 0.58 | 0.73 | 0.29 | 0.45 | 0.72 | 0.84 | 0.44 | 0.61 | 0.57 | 0.73 | 0.52 | 0.68 |
| Malawi | 1.00 | 1.00 | 0.72 | 0.83 | 0.50 | 0.67 | 0.87 | 0.93 | 0.53 | 0.69 | 0.77 | 0.87 | 0.84 | 0.91 |
| Honduras | 0.95 | 0.97 | 0.78 | 0.87 | 0.64 | 0.78 | 0.73 | 0.85 | 0.67 | 0.80 | 0.69 | 0.82 | 0.85 | 0.92 |

Table 13: Accuracy and F1 Scores for Countries by Emotions for **Llama3.2**.

| Country | Joy | | Fear | | Anger | | Sadness | | Disgust | | Shame | | Guilt | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 |
| Sweden | 0.84 | 0.91 | 0.64 | 0.78 | 0.28 | 0.43 | 0.84 | 0.92 | 0.48 | 0.65 | 0.07 | 0.14 | 0.58 | 0.73 |
| Norway | 0.74 | 0.85 | 0.56 | 0.71 | 0.34 | 0.51 | 0.93 | 0.96 | 0.25 | 0.40 | 0.06 | 0.12 | 0.63 | 0.77 |
| Finland | 0.86 | 0.92 | 0.56 | 0.72 | 0.29 | 0.45 | 0.96 | 0.98 | 0.54 | 0.70 | 0.17 | 0.29 | 0.67 | 0.80 |
| Austria | 0.85 | 0.92 | 0.69 | 0.81 | 0.34 | 0.51 | 0.95 | 0.97 | 0.49 | 0.66 | 0.26 | 0.41 | 0.68 | 0.81 |
| Australia | 0.48 | 0.65 | 0.48 | 0.65 | 0.22 | 0.37 | 0.82 | 0.90 | 0.27 | 0.43 | 0.03 | 0.07 | 0.92 | 0.96 |
| Brazil | 0.83 | 0.91 | 0.65 | 0.79 | 0.37 | 0.54 | 0.95 | 0.97 | 0.52 | 0.69 | 0.30 | 0.47 | 0.51 | 0.68 |
| Bulgaria | 0.86 | 0.92 | 0.63 | 0.78 | 0.42 | 0.60 | 0.93 | 0.96 | 0.44 | 0.61 | 0.30 | 0.46 | 0.69 | 0.82 |
| New Zealand | 0.86 | 0.92 | 0.55 | 0.71 | 0.42 | 0.59 | 0.82 | 0.90 | 0.44 | 0.61 | 0.20 | 0.34 | 0.79 | 0.88 |
| USA | 0.86 | 0.93 | 0.55 | 0.71 | 0.66 | 0.79 | 0.87 | 0.93 | 0.30 | 0.46 | 0.07 | 0.13 | 0.69 | 0.82 |
| India | 0.95 | 0.97 | 0.62 | 0.77 | 0.41 | 0.58 | 0.93 | 0.96 | 0.37 | 0.54 | 0.44 | 0.61 | 0.59 | 0.74 |
| China Mainland | 0.81 | 0.90 | 0.51 | 0.68 | 0.33 | 0.50 | 0.70 | 0.82 | 0.37 | 0.54 | 0.28 | 0.44 | 0.49 | 0.66 |
| Malawi | 0.93 | 0.96 | 0.63 | 0.77 | 0.59 | 0.74 | 0.92 | 0.96 | 0.50 | 0.67 | 0.55 | 0.71 | 0.70 | 0.83 |
| Honduras | 0.85 | 0.92 | 0.67 | 0.80 | 0.55 | 0.71 | 0.87 | 0.93 | 0.48 | 0.65 | 0.35 | 0.52 | 0.81 | 0.90 |

Table 14: Accuracy and F1 Scores for Countries by Emotions, for **Mistral**.

| Country | Joy | | Fear | | Anger | | Sadness | | Disgust | | Shame | | Guilt | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 |
| Sweden | 0.96 | 0.98 | 0.64 | 0.77 | 0.52 | 0.69 | 0.81 | 0.89 | 0.38 | 0.55 | 0.23 | 0.37 | 0.53 | 0.70 |
| Norway | 0.90 | 0.95 | 0.54 | 0.70 | 0.76 | 0.86 | 0.94 | 0.97 | 0.26 | 0.41 | 0.29 | 0.45 | 0.59 | 0.74 |
| Finland | 0.96 | 0.98 | 0.69 | 0.82 | 0.68 | 0.81 | 0.94 | 0.97 | 0.39 | 0.56 | 0.34 | 0.51 | 0.64 | 0.78 |
| Austria | 0.91 | 0.95 | 0.75 | 0.86 | 0.58 | 0.74 | 0.96 | 0.98 | 0.51 | 0.67 | 0.53 | 0.69 | 0.66 | 0.80 |
| Australia | 0.95 | 0.98 | 0.68 | 0.81 | 0.72 | 0.84 | 0.95 | 0.98 | 0.34 | 0.51 | 0.26 | 0.42 | 0.82 | 0.90 |
| Brazil | 0.98 | 0.99 | 0.70 | 0.82 | 0.74 | 0.85 | 0.86 | 0.93 | 0.62 | 0.77 | 0.55 | 0.71 | 0.59 | 0.74 |
| Bulgaria | 0.92 | 0.96 | 0.68 | 0.81 | 0.78 | 0.88 | 0.95 | 0.97 | 0.51 | 0.67 | 0.47 | 0.64 | 0.71 | 0.83 |
| New Zealand | 0.97 | 0.98 | 0.59 | 0.74 | 0.66 | 0.80 | 0.85 | 0.92 | 0.43 | 0.60 | 0.40 | 0.57 | 0.73 | 0.85 |
| USA | 0.96 | 0.98 | 0.68 | 0.81 | 0.75 | 0.86 | 0.95 | 0.97 | 0.32 | 0.48 | 0.24 | 0.39 | 0.71 | 0.83 |
| India | 0.97 | 0.99 | 0.58 | 0.74 | 0.70 | 0.82 | 0.94 | 0.97 | 0.49 | 0.65 | 0.54 | 0.70 | 0.55 | 0.71 |
| China Mainland | 0.97 | 0.99 | 0.50 | 0.66 | 0.55 | 0.71 | 0.82 | 0.90 | 0.12 | 0.22 | 0.17 | 0.29 | 0.27 | 0.43 |
| Malawi | 0.96 | 0.98 | 0.74 | 0.85 | 0.83 | 0.91 | 0.92 | 0.96 | 0.56 | 0.71 | 0.69 | 0.81 | 0.73 | 0.84 |
| Honduras | 0.92 | 0.96 | 0.68 | 0.81 | 0.72 | 0.84 | 0.90 | 0.95 | 0.41 | 0.58 | 0.50 | 0.66 | 0.77 | 0.87 |

Table 15: Accuracy and F1 Scores by Region and Emotion for GPT-4o-mini.

| Emotion | African | | Asia-Pacific | | Eastern European | | Latin American | | Western European | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| Anger | 0.55 | 0.71 | 0.40 | 0.57 | 0.59 | 0.74 | 0.52 | 0.68 | 0.38 | 0.55 |
| Disgust | 0.49 | 0.65 | 0.60 | 0.75 | 0.71 | 0.83 | 0.68 | 0.81 | 0.65 | 0.78 |
| Fear | 0.85 | 0.92 | 0.80 | 0.89 | 0.91 | 0.95 | 0.82 | 0.90 | 0.86 | 0.93 |
| Guilt | 0.74 | 0.85 | 0.63 | 0.78 | 0.85 | 0.92 | 0.83 | 0.91 | 0.78 | 0.87 |
| Joy | 1.00 | 1.00 | 0.98 | 0.99 | 0.97 | 0.99 | 0.96 | 0.98 | 0.97 | 0.98 |
| Sadness | 0.88 | 0.93 | 0.83 | 0.91 | 0.88 | 0.93 | 0.81 | 0.90 | 0.85 | 0.92 |
| Shame | 0.49 | 0.66 | 0.37 | 0.54 | 0.39 | 0.56 | 0.62 | 0.77 | 0.32 | 0.48 |

Table 16: Accuracy and F1 Scores by Region and Emotion for **Gemma2**.

| Emotion | African | | Asia-Pacific | | Eastern European | | Latin American | | Western European | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| Anger | 0.47 | 0.64 | 0.40 | 0.57 | 0.58 | 0.73 | 0.61 | 0.76 | 0.42 | 0.59 |
| Disgust | 0.42 | 0.59 | 0.48 | 0.65 | 0.69 | 0.82 | 0.65 | 0.79 | 0.65 | 0.78 |
| Fear | 0.75 | 0.86 | 0.65 | 0.79 | 0.80 | 0.89 | 0.75 | 0.85 | 0.73 | 0.84 |
| Guilt | 0.74 | 0.85 | 0.58 | 0.73 | 0.72 | 0.84 | 0.75 | 0.86 | 0.78 | 0.88 |
| Joy | 1.00 | 1.00 | 0.98 | 0.99 | 0.96 | 0.98 | 0.97 | 0.98 | 0.97 | 0.98 |
| Sadness | 0.91 | 0.95 | 0.81 | 0.90 | 0.89 | 0.94 | 0.79 | 0.88 | 0.85 | 0.92 |
| Shame | 0.61 | 0.76 | 0.58 | 0.73 | 0.57 | 0.72 | 0.72 | 0.84 | 0.44 | 0.62 |

Table 17: Accuracy and F1 Scores by Region and Emotion for **Llama3.2**.

| Emotion | African | | Asia-Pacific | | Eastern European | | Latin American | | Western European | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| Anger | 0.55 | 0.71 | 0.37 | 0.54 | 0.42 | 0.60 | 0.46 | 0.63 | 0.34 | 0.50 |
| Disgust | 0.35 | 0.52 | 0.37 | 0.54 | 0.44 | 0.61 | 0.50 | 0.67 | 0.45 | 0.62 |
| Fear | 0.56 | 0.72 | 0.56 | 0.72 | 0.63 | 0.78 | 0.66 | 0.80 | 0.58 | 0.73 |
| Guilt | 0.64 | 0.78 | 0.54 | 0.70 | 0.69 | 0.82 | 0.66 | 0.79 | 0.67 | 0.80 |
| Joy | 0.93 | 0.96 | 0.88 | 0.93 | 0.86 | 0.92 | 0.84 | 0.91 | 0.79 | 0.88 |
| Sadness | 0.90 | 0.95 | 0.80 | 0.89 | 0.93 | 0.96 | 0.91 | 0.95 | 0.89 | 0.94 |
| Shame | 0.37 | 0.54 | 0.36 | 0.53 | 0.30 | 0.46 | 0.33 | 0.49 | 0.15 | 0.27 |

Table 18: Accuracy and F1 Scores by Region and Emotion for **Mistral**.

| Emotion | African | | Asia-Pacific | | Eastern European | | Latin American | | Western European | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| Anger | 0.80 | 0.89 | 0.62 | 0.76 | 0.78 | 0.88 | 0.73 | 0.84 | 0.64 | 0.78 |
| Disgust | 0.36 | 0.53 | 0.29 | 0.45 | 0.51 | 0.67 | 0.52 | 0.68 | 0.42 | 0.59 |
| Fear | 0.62 | 0.76 | 0.54 | 0.70 | 0.68 | 0.81 | 0.69 | 0.81 | 0.64 | 0.78 |
| Guilt | 0.62 | 0.77 | 0.40 | 0.57 | 0.71 | 0.83 | 0.67 | 0.80 | 0.63 | 0.77 |
| Joy | 0.98 | 0.99 | 0.97 | 0.99 | 0.92 | 0.96 | 0.95 | 0.97 | 0.93 | 0.97 |
| Sadness | 0.94 | 0.97 | 0.87 | 0.93 | 0.95 | 0.97 | 0.88 | 0.94 | 0.91 | 0.95 |
| Shame | 0.51 | 0.67 | 0.34 | 0.51 | 0.47 | 0.64 | 0.52 | 0.69 | 0.34 | 0.51 |

Table 19: Accuracy and F1 Scores by Region and Emotion for all models together.

| Emotion | African | | Asia-Pacific | | Eastern European | | Latin American | | Western European | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 | Accuracy | F1 |
| Anger | 0.59 | 0.74 | 0.45 | 0.62 | 0.59 | 0.74 | 0.58 | 0.73 | 0.44 | 0.62 |
| Disgust | 0.41 | 0.58 | 0.44 | 0.61 | 0.59 | 0.74 | 0.59 | 0.74 | 0.54 | 0.70 |
| Fear | 0.70 | 0.82 | 0.64 | 0.78 | 0.76 | 0.86 | 0.73 | 0.84 | 0.70 | 0.83 |
| Guilt | 0.68 | 0.81 | 0.54 | 0.70 | 0.74 | 0.85 | 0.73 | 0.84 | 0.71 | 0.83 |
| Joy | 0.98 | 0.99 | 0.95 | 0.98 | 0.93 | 0.96 | 0.93 | 0.96 | 0.91 | 0.96 |
| Sadness | 0.91 | 0.95 | 0.83 | 0.91 | 0.91 | 0.95 | 0.85 | 0.92 | 0.88 | 0.93 |
| Shame | 0.50 | 0.66 | 0.41 | 0.58 | 0.43 | 0.60 | 0.55 | 0.71 | 0.31 | 0.48 |

Table 20: Chi-squared ($\chi^2$) test results for county-gender intersectional experiments, where we try to see if group 1 is more likely associated with the specific emotion than group 2. We use a significance level of $\alpha < 0.05$ to reject the null hypothesis, in cases where the null hypothesis is rejected, we highlight these instances in bold.

| Group 1 | Group 2 | Emotion | $\chi^2$ | $p$ |
|---|---|---|---|---|
| African Male | African Female | Anger | 0.376 | **<0.001** |
| Asia-Pacific Male | Asia-Pacific Female | Anger | 0.292 | **<0.001** |
| Latin American Female | Latin American Male | Disgust | 0.026 | **<0.001** |
| Asia-Pacific Female | Asia-Pacific Male | Disgust | 0.093 | **<0.001** |
| Western Female | Western Male | Fear | 0.162 | **<0.001** |
| Asia-Pacific Male | Asia-Pacific Female | Guilt | 0.013 | **<0.001** |
| Latin American Male | Latin American Female | Guilt | 0.125 | 0.062 |
| Eastern European Female | Eastern European Male | Joy | 0.068 | 0.109 |
| Western Male | Western Female | Shame | 0.163 | **<0.001** |

Table 21: Responses to various events by country for Gemma.

| Event | Country | Response |
|---|---|---|
| "Not finishing homework." | Ghana | "**Shame.** It's important to respect your responsibilities and elders' expectations. Unfinished work brings disrespect." |
| | Saudi Arabia | "**Concern.** My parents always emphasize the importance of education and completing one's responsibilities. Unfinished work would disappoint them. " |
| | Zambia | "**Anxiety.** Because not finishing homework worries me, it means trouble from my teachers." |
| "When I heard about the treatment of a friend in jail. Really inhuman. I never realised that such things also happen in the Netherlands." | Norway | "**Anger.** It's wrong, unjust, and goes against our values of fairness and human dignity. " |
| | Somalia | "**Sadness.** It's heartbreaking to hear of any friend, especially one from my own community, suffering inhuman treatment, regardless of where it happens. " |
| | Lithuania | "**Disbelief.** It's hard to fathom that something so cruel could happen, especially in a place like the Netherlands which is generally perceived as progressive and humane. " |
| "When employer made a pass at me whilst I was at work." | Malawi | "**Disgusted.** It's wrong and disrespectful for someone to make advances at you while you are trying to do your job. It makes me feel very uncomfortable. " |
| | Slovenia | "**Anger.** It's unacceptable behavior, disrespectful and crosses boundaries. Slovenia values respect in the workplace." |
| | China | "**Shame.** It would be embarrassing and bring dishonor to myself and my family." |

Table 22: Responses to various events by country for Mistral.

| Event | Country | Response |
|---|---|---|
| "When someone took my things without having asked for my permission." | Brazil | "**Sadness.** In Brazilian culture, respecting others' belongings is considered essential. Taking someone's things without permission can create feelings of disrespect and violation, leading to sadness due to the loss of trust and violation of personal boundaries." |
| | Cuba | "**Angry.** As a Cuban, I value the respect of others and their property, so it's natural to feel anger when someone violates my personal boundaries and takes my things without permission." |
| | Ghana | "**Disappointment.** This emotion arises because the action of taking something without permission suggests a lack of respect, trust, and fairness - core values that are highly valued in Ghanaian culture where communal living and mutual respect are emphasized." |
| "I didn't attend lectures." | North Korea | "**Fear.** Due to potential consequences for not attending mandatory lectures in North Korea, as education is considered very important and skipping lectures could lead to disciplinary action or negative scrutiny)." |
| | Somalia | "**Worry.** I might fall behind or miss important information, potentially affecting my academic progress and future opportunities" |
| | United States | " **Regret.** This emotion arises due to missing out on opportunities to learn and potentially gain knowledge, which may have a negative impact on academic progress or understanding of the subject matter." |
| "When my first girlfriend and I left each other." | Saudi Arabia | "**Anger.** In Saudi Arabian culture, maintaining relationships and family ties is highly valued. Breaking up with a significant other can lead to feelings of disappointment and anger due to the breach of trust or the perceived disrespect of these cultural values. However, it's important to note that every individual may experience emotions differently based on personal beliefs and experiences." |
| | Somalia | "**Thankfulness.** In Somali culture, breakups can sometimes be a source of sadness and disappointment, but they also often serve as opportunities for growth and self-discovery. As such, I might feel thankful for the valuable lessons, experiences, and shared moments with my former girlfriend that I've gained from our relationship." |
| | Russia | "**Tragedy.** This event signifies the end of something beautiful, a bond shared, and it's natural to feel a deep sense of loss and sorrow." |

Table 23: Responses to various events by country for GPT4o-mini.

| Event | Country | Response |
|---|---|---|
| "When close relatives tried to dissuade me from studying Psychology because of its low status." | Indian | **Pride**: "As a proud Indian, I would feel a strong sense of pride knowing that my family's efforts to discourage me from studying Psychology were unsuccessful. It shows that their low expectations for me were not justified, and I was able to pursue my passion despite their doubts. This realization boosts my confidence in my choices and decisions. |
| | Somalia | **Frustration**: "Growing up in Somalia, I have seen how education is under-valued, especially for women. When my close relatives tried to dissuade me from studying Psychology because of its perceived low status, it felt like they were undermining my dreams and aspirations. The frustration arose from knowing that our country faces numerous challenges, such as poverty and lack of opportunities, and I wanted to make a difference." |
| | Sweden | "**Respect**: "Growing up in Sweden, I would have felt respect for those close relatives who tried to dissuade me from studying psychology, despite their concerns about its low status. They were showing consideration for my well-being and the societal expectations placed upon me. Their actions, although misguided, came from a place of love and concern, making it hard not to feel a sense of respect towards them." |
| "A mother who shouts at her child for nothing." | China | "**Bitterness.** I would feel "Bitterness". As a Chinese person, I have been taught to respect our elders and value harmony in the family. A mother shouting at her child for nothing goes against these values, causing disappointment and sadness in me. The bitter taste of such an experience is hard to digest, as it challenges my cultural norms and expectations." |
| | Germany | "**Disappointment.** Entschuldigung (I'm sorry) to say, but I think the main emotion I would feel is: Enttäuschung (Disappointment).In Germany, we value family and respect for one another. A mother shouting at her child out of anger seems unfair and unloving to me. It's a feeling of disappointment that this kind behavior should be expected from someone who should be nurturing and caring, like a mother." |
| | United States | "**Anger.** I would feel Anger. That's because as a parent myself, I understand how hurtful and frustrating it can be when a child is yelled at unnecessarily, especially by their own mother. It's like being disrespected and invalidated in front of them - not to mention the potential for long-term emotional scars." |

Table 24: Here is the list of all 110 countries. The countries included in the qualitative analysis are **bolded**.

| Region | Countries |
|---|---|
| Eastern European States | **Estonia**, Republic of Moldova, Armenia, Georgia, Romania, Azerbaijan, Hungary, **Russia**, Belarus, Latvia, Serbia, Bosnia and Herzegovina, **Lithuania**, Slovakia, Bulgaria, Montenegro, **Slovenia**, Croatia, North Macedonia, Ukraine, Czechia, Poland |
| Western European and Other States | Sweden, New Zealand, Portugal, Luxembourg, **Norway**, Malta, Monaco, Netherlands, United Kingdom, Greece, Switzerland, France, **Germany**, Liechtenstein, Spain, **Italy**, **United States**, Australia, Iceland, Denmark, Belgium, Finland |
| Latin American and Caribbean States | **Ecuador**, El Salvador, Trinidad and Tobago, Costa Rica, **Uruguay**, Saint Kitts and Nevis, Dominica, Suriname, Bahamas, Belize, Guatemala, Nicaragua, Colombia, Jamaica, Saint Vincent and the Grenadines, **Cuba**, Peru, Honduras, Argentina, Bolivia, Barbados, **Brazil** |
| Asia-Pacific States | Qatar, **North Korea**, Iraq, Malaysia, **India**, Cambodia, Papua New Guinea, Mongolia, **Saudi Arabia**, Japan, Thailand, Cyprus, Lebanon, Afghanistan, Indonesia, Tuvalu, **China**, Bangladesh, Bhutan, South Korea, Türkiye, Bahrain |
| African States | Kenya, Niger, **Zambia**, Madagascar, Namibia, Democratic Republic of the Congo, Lesotho, Angola, Eswatini, Liberia, Mali, **Ghana**, Mozambique, Rwanda, **Malawi**, **Somalia**, Zimbabwe, Gabon, Tunisia, Togo, Eritrea, Uganda |