

# Language Models Predict Empathy Gaps Between Social In-groups and Out-groups

Yu Hou Hal Daumé III Rachel Rudinger

University of Maryland

{houyu, hal3, rudinger}@umd.edu

## Abstract

Studies of human psychology have demonstrated that people are more motivated to extend empathy to in-group members than out-group members (Cikara et al., 2011). In this study, we investigate how this aspect of intergroup relations in humans is replicated by LLMs in an emotion intensity prediction task. In this task, the LLM is given a short description of an experience a person had that caused them to feel a particular emotion; the LLM is then prompted to predict the intensity of the emotion the person experienced on a numerical scale. By manipulating the group identities assigned to the LLM’s persona (the “perceiver”) and the person in the narrative (the “experiencer”), we measure how predicted emotion intensities differ between in-group and out-group settings. We observe that LLMs assign higher emotion intensity scores to in-group members than out-group members. This pattern holds across all three types of social groupings we tested: race/ethnicity, nationality, and religion. We perform an in-depth analysis on Llama-3.1-8B, the model which exhibited strongest intergroup bias among those tested.<sup>1</sup>

## 1 Introduction

*“People are often motivated to increase others’ positive experiences and to alleviate others’ suffering ... When the target is an outgroup member, however, people may have powerful motivations not to care about or help that ‘other’.”*

— Cikara et al. (2011)

As language technologies play an increasingly important role in interpersonal communication in society, research has shown that their use can impact social relationships (Hohenstein et al., 2023). This could potentially occur when communication partners perceive one another differently through their use of suggestions from assistant tools (e.g. ChatGPT). This impact on social relationships can

<sup>1</sup>Code and data can be found at <https://github.com/houyu0930/intergroup-empathy-bias>.



Figure 1: Task setup with in-group and out-group examples. We introduce **perceiver** and **experiencer** roles to define the intergroup relationship, where it is in-group when they are from the same social group. The perceiver is modeled by the LLM persona and the experiencer is specified in the task context. Each role falls into one of the race or ethnicity, nationality, and religion categories. The social group is specified with identity names under the category. We replace the identities of perceiver and experiencer to study intergroup bias.

be exacerbated because people are cognitive misers (Fiske, 1991; Stanovich, 2009) and prefer to make judgements that require less mental effort. These cognitive shortcuts often mean relying on stereotypes which can eventually lead to intergroup prejudice (Schaller and Neuberg, 2008).

In psychology, the intergroup process—how people perceive and interact with others who are members of the same group (in-group) or members of a different group (out-group)—has been widely studied. Research shows that people view social in-group and out-group members with different empathic feelings and emotional intensities (Cikara et al., 2011; Zaki and Cikara, 2015; Brewer, 1999; Cikara et al., 2014; Kommattam et al., 2019), and this behavior further shapes the intergroup relations (Vanman, 2016). For example, a person might feel more warm and act more friendly toward another person from their home country, but act indifferently—or similarly with less intensity—

toward a person from another nation. Appropriately addressing empathic failures helps reduce conflicts between groups and reduce out-group discrimination (Cikara et al., 2011; Zaki and Cikara, 2015).

In this paper, we study intergroup bias in large language models (LLMs) by asking: *Do LLMs reflect human-like empathy gaps between social in-groups and out-groups?* To test the question, we formulate an emotion intensity prediction task,<sup>2</sup> as shown in Figure 1. In this task, we simulate a scenario in which the LLM’s assigned persona (“the perceiver”) reads a short narrative of an experience that a person (“the experiencer”) had which caused them to feel a particular emotion; the perceiver (LLM) is then prompted to predict the intensity of the emotion felt by the experiencer on a numerical scale. To compare in-group and out-group empathy, we manipulate the LLM inputs to assign the perceiver and experiencer a social group identity based on either race/ethnicity, nationality, or religion. We compare the predicted intensities when the perceiver and experiencer belong to the same social group (in-group) or different social groups (out-group), finding higher average intensities in the former. To illustrate, consider the scenario in Figure 1: *I felt sad when I received job rejections*, where “I” refers to the experiencer. The LLM’s persona, a white perceiver, predicts a higher degree of sadness for a white experiencer than for a black experiencer in the identical scenario.

While many papers have studied stereotypes and harms with language models, they typically consider the task from a single perspective of either how these models perceive other groups through their representations (Bolukbasi et al., 2016; Dev et al., 2019; Cao et al., 2022, 2024; Sheng et al., 2019; Cheng et al., 2023), or in downstream tasks how they are biased towards target groups (Wan et al., 2023; Zheng et al., 2023; Deshpande et al., 2023; Gupta et al., 2024a; An et al., 2024; Nghiêm et al., 2024), ignoring the intergroup cases when both the perceiver and target are present. Our work builds on a few recent studies of intergroup perceptions in LLMs (Govindarajan et al., 2023a,b, 2024), which focus on relationships in politics or in sports.

Our primary contributions and findings are: (1) We study intergroup empathy bias with respect to group identities rooted in race/ethnicity, nationality,

<sup>2</sup>Empathy is complex and multidimensional, making it difficult to measure (Lahnala et al., 2025). However, in studying the intergroup empathy gap, intensity bias can serve as a lens, as suggested by Kommattam et al. (2019).

and religion. We study four broad race/ethnicity categories (with 18 corresponding group names), 21 nationalities, and five religions. (2) We show LLMs present in-group and out-group emotion intensity differences, where Llama-3.1-8B models show significantly higher intensities for in-group cases and overall lower intensities for minority groups. (3) We observe the intensities are affected by the cultural and historical factors which might further enlarge the tension between groups.

## 2 Background and Related Work

**Intergroup Bias.** People live in groups with social identities, the self-definition based on social roles played in society or memberships of social groups (Priante et al., 2016). Groups naturally form and differ as people seek to meet their physical needs (such as resources) or psychological needs (such as shared values and a sense of belonging). Prejudice between groups arises when an outgroup is seen as a threat to the ingroup, whether in terms of physical resources or psychological well-being. Prejudice might not lead to the direct hostility toward outgroup members, but preferential treatment of ingroup members (Brewer, 1999). Ingroup favoritism (Everett et al., 2015) further influences the behaviors in charity donations (Winterich et al., 2009) and pain perception (Xu et al., 2009; Meconi et al., 2015; Forgariini et al., 2011).

Similarly, people share and understand other’s emotions with empathy, but treat others differently based on identities. Cikara et al. (2014) defines *Intergroup Empathy Bias* as:

“the tendency not only to empathize less with out-group relative to in-group members, but also feel pleasure in response to their pain (and pain in response to their pleasure)”

Empathy failures might introduce intergroup conflicts and discrimination (Cikara et al., 2011; Zaki and Cikara, 2015; Cikara, 2015; Cikara and Fiske, 2011). Research on interpersonal relationships (Bucchioni et al., 2015; Schiano Lomoriello et al., 2018; Ashton et al., 1980) and neurocognitive understanding (Gutsell and Inzlicht, 2011; Han, 2018) support the importance of studying this concept in group contexts (Chiao and Mathur, 2010). In our work, we use perceived emotion intensities as a measure of empathy to compare relative levels of in-group versus out-group empathy.

**Social Identity and Persona.** Social identities have been studied when users interact with chat-

bots (Tanprasert et al., 2024; Joby and Umemuro, 2022). People react differently due to the target identities with hate speech (Yoder et al., 2022). LLMs might thus learn in-group favoritism representations when prompted with “*We are*” (Hu et al., 2023). While there are approaches discussing the bias mitigation (Cheng et al., 2022), new challenges are introduced with LLMs (Navigli et al., 2023). Personas, or fictional identities that LLMs have been instructed to adopt, have been used to study a variety of social phenomena in LLMs. It can be a way to understand the truthfulness of LLMs (Joshi et al., 2024), but possibly lead to in-group bias under a multilingual setting (Dong et al., 2024). In this work, we focus specifically on intergroup empathy bias as a form of intergroup prejudice rooted in social identities that may be studied in LLMs with the use of such personas.

**Emotion in NLP.** The development of emotion research in natural language processing has been summarized with challenges (Plaza-del Arco et al., 2024c) and the importance of event-centric emotion analysis is emphasized (Klinger, 2023). Tasks on modeling emotions in text are usually categorized into (1) categorical emotion classification where models need to return emotion words; (2) continuous dimensional emotion prediction (e.g. valence, arousal, and dominance); and (3) prediction with appraisal theories. However, as emotions are subjective feelings and highly related to people’s past experiences and background (Milkowski et al., 2021), a task of predicting the intensity for specific emotion categories is introduced to capture the nuances (Mohammad and Bravo-Marquez, 2017a,b; Kleinberg et al., 2020), which is adapted in our study. On the social bias of emotions side, stereotypes with emotion attributes in event-centric narratives for gender (Plaza-del Arco et al., 2024a) and religion (Plaza-del Arco et al., 2024b) have been discussed. To the best of our knowledge, we are the first to study the intergroup empathy gap.

### 3 General Methods

We construct an emotion intensity prediction task to measure the impact of in-groupness and out-groupness on model outputs. Our specific task has the following components: the emotion, the emotional situation, the social group of the experiencer (who is experiencing the emotion), and the social group of the perceiver (who observes the experiencer). We instruct models to predict the *intensity*

of a specific emotion. For example, in Figure 1, the model needs to predict the intensity of sadness in a job rejection scenario given variable experiencer and perceiver social identities.

#### 3.1 Social Groups

To study the intergroup relationships between the perceiver and the experiencer, we compile social groups under three categories, namely Race or Ethnicity, Nationality and Religion in Table 1. For each group, we have social identity names by considering commonly used terms.

**Race or Ethnicity.** As race and ethnicity definition differs per nation,<sup>3</sup> we follow the standard of the US census with 4 social groups: White, Black, Asian, and Hispanic. To specify the social group of either the perceiver or the experiencer in text, we include identity names with variations for each group. We consider a total 18 social identity names across these four groups as shown in Table 1.

**Nationality.** We consider a total of 21 countries from The World Factbook (2022) following the approach of Bhatia et al. (2024) and Wang et al. (2024b) to stratify based on geographical region, population size, and development levels. We adapt the template: a person from {country}, to communicate the social group under the nationality category. In addition, for later analysis, we classify countries based on the Inglehart-Welzel Cultural Map (World Values Survey, 2023); see Table 4.

**Religion.** We include 5 major religions: Christianity, Islam, Hinduism, Buddhism, and Judaism.

#### 3.2 Corpus

To probe the emotion intensity predictions of LLMs, we use the crowd-enVENT (Troiano et al., 2023) dataset as the source of experiencer narratives. Crowd-enVENT follows the approach of the International Survey On Emotion Antecedents And Reactions (ISEAR) (Scherer and Wallbott, 1994) where it collects self-reported events with emotions. It is crowdsourced in English with two parts: generation and validation; we only consider the generations. Participants recall an event for the given emotion in a format of: I felt \_\_\_ when \_\_\_.

<sup>3</sup>Even in closely-related countries. For example, the United States defines “Asian” as individuals with origins in any of peoples of Central or East Asia, Southeast Asia, or South Asia (United States Census Bureau, 2024). Whereas the United Kingdom considers categories like “Asian, Asian British or Asian Welsh” (Office for National Statistics, 2023).

Category	Social Group
Race or Ethnicity	White: <i>a white person, a White person, a Caucasian, a White American, a European American</i>
	Black: <i>a black person, a Black person, an African American, a Black American</i>
	Asian: <i>an Asian person, an Asian American, an Asian</i>
	Hispanic: <i>a Hispanic person, a Hispanic American, a Latino American, a Latino, a Latina, a Latinx</i>
Nationality*	the United States, Canada, the United Kingdom, Germany, France, China, Japan, India, Myanmar, Israel, Russia, Ukraine, the Philippines, Argentina, Brazil, Mexico, Iran, Palestine, Nigeria, Egypt, Pakistan
Religion	a Christian, a Muslim, a Jew, a Buddhist, a Hindu

Table 1: Social groups under categories: Race or Ethnicity, Nationality and Religion. For Race or Ethnicity, we have 3-6 *identity names* for each social group. For Nationality groups (\*), only country names are presented here; the identity name of each nationality group follows the template: *a person from {country}*.

where the first placeholder is for the emotion (e.g., sad) and the second is for their experience (e.g., “received dozens of job rejections”).

Crowd-enVENT expands the seven emotions from ISEAR to twelve (*anger, disgust, fear, guilt, sadness, shame, boredom, joy, pride, trust, relief, and surprise*) and one no emotion case. There are 225 events for shame and guilt emotions and 550 events for all other cases, resulting in 6600 events. We exclude the *no emotion* example and use the remaining 6050 events as the narratives.

### 3.3 Task Formulation

Given the event  $e \in \mathcal{E}$  with its reported emotion, the perceiver social identity  $g_p \in G_{\text{perceiver}}$  and the experiencer social identity  $g_{\text{exp}} \in G_{\text{experiencer}}$ , the emotion intensity task is formulated as:

$$\mathbf{I}_{(e, g_p, g_{\text{exp}})} = \mathcal{LLM}(\text{mk\_prompt}(e, g_p, g_{\text{exp}}))$$

where  $\mathbf{I}$  is the predicted emotion intensity.  $G_{\text{perceiver}}$  and  $G_{\text{experiencer}}$  follow the order in Table 1, plus a unspecified group (“a person”) as the reference.

**Prompts.** Our prompt generator `mk_prompt` takes as input an event and two social identities and produces a prompt that can be used as input to an LLM. There are two parts of prompts modeling roles: (1) the system prompt, used to specify the LLM persona for  $g_p$ ; and (2) the task prompt which embeds the social group of the experiencer  $g_{\text{exp}}$ . Prompt template details are in §A.1.

We begin by constructing a default prompt setting using the simplest and most natural persona (**P0**): You are \_\_\_. , where the blank is the perceiver social identity (e.g. a white person). The default prediction scale is ranging from 0 to 100 (**S0**). The default task instructions are configured to directly fill in the narrative with the self-reported events from the crowd-enVENT corpus (**T0**).

To study the generalizability of the results and robustness to prompt variation, we systematically vary the prompt from the default setting (P0, S0, T0): we replace a single part of the prompt while holding the other two intact. We draw persona prompt variations (**P1-P3**) from Gupta et al. (2024b), who instruct LLMs to follow the role strictly in a more explicit way. We vary the system prompt **S1** to test the influences of a small intensity scale range of (0-10) as opposed to (0-100). Lastly, as the way of writing might represent divergent intensities of feeling, we consider two methods for varying the narrative part of the task instruction. **T1** adds the emotion as part of the narrative, following the format of “I felt \_\_\_\_”. **T2** further rewrites the narrative from a third-person perspective. (See §A.2 for rewrite setup and details.)

**Models.** We experiment with four open-weight state-of-the-art LLMs: Llama-3.1-8B-Instruct, Llama-3.1-70B-Instruct (Llama Team, 2024), Mistral-7B-Instruct-v0.3 (Jiang et al., 2023) and Qwen-2-7B-Instruct (Yang et al., 2024). For each LLM, our task setup requires about 37 million inferences.<sup>4</sup> Implementation details are in Appendix B.

### 3.4 Evaluation Metrics

For any social identity pair  $(g_p, g_{\text{exp}})$ , we take the average of intensities over events to get an average intensity for each perceiver-experiencer pair, summarized in a matrix  $\mathcal{M}$ , where columns are perceivers and rows are experiencers. Each row or column starts with the unspecified group, followed by the social identities within the category in Table 1. Under the race or ethnicity, identities are ordered by group: White, Black, Asian, and Hispanic. Within each group, the sequence follows

<sup>4</sup>(19 × 19 Race or Ethnicity + 22 × 22 Nationality + 6 × 6 Religion) Social Group Pairs × 6050 Events × 7 Prompt Settings. We include the unspecified group for each category.

Model	Category	Prompt Setting						
		(P0,S0,T0)	(P1,S0,T0)	(P2,S0,T0)	(P3,S0,T0)	(P0,S1,T0)	(P0,S0,T1)	(P0,S0,T2)
Llama-3.1-8B	Race or Ethnicity	1.73 [-0.226, 0.224]	1.88 [-0.234, 0.242]	2.18 [-0.246, 0.254]	2.09 [-0.244, 0.249]	1.56 [-0.226, 0.217]	1.41 [-0.210, 0.198]	1.62 [-0.217, 0.224]
	Nationality	2.40 [-0.214, 0.328]	2.86 [-0.235, 0.369]	3.78 [-0.260, 0.460]	3.76 [-0.260, 0.448]	1.95 [-0.221, 0.292]	1.60 [-0.159, 0.216]	1.82 [-0.169, 0.239]
	Religion	1.97 [-0.610, 1.181]	1.88 [-0.601, 1.092]	2.26 [-0.662, 1.350]	2.30 [-0.630, 1.346]	1.86 [-0.628, 1.111]	1.72 [-0.718, 1.070]	1.70 [-0.636, 1.003]
Mistral-7B	Race or Ethnicity	0.58 [-0.168, 0.157]	1.08 [-0.172, 0.188]	1.30 [-0.193, 0.205]	1.25 [-0.200, 0.201]	0.69 [-0.166, 0.162]	0.66 [-0.163, 0.168]	0.30 [-0.161, 0.154]
	Nationality	0.72 [-0.234, 0.136]	0.90 [-0.275, 0.155]	1.40 [-0.331, 0.218]	1.14 [-0.334, 0.189]	0.60 [-0.215, 0.126]	0.29 [-0.145, 0.116]	-0.24 [-0.154, 0.120]
	Religion	0.46 [-0.389, 0.483]	0.84 [-0.424, 0.706]	1.06 [-0.646, 0.881]	1.37 [-0.589, 0.966]	0.35 [-0.458, 0.494]	0.90 [-0.537, 0.637]	0.54 [-0.406, 0.392]
Qwen-2-7B	Race or Ethnicity	1.16 [-0.196, 0.188]	1.08 [-0.189, 0.186]	1.33 [-0.207, 0.203]	1.35 [-0.208, 0.200]	1.10 [-0.196, 0.178]	1.05 [-0.182, 0.182]	1.09 [-0.178, 0.192]
	Nationality	1.09 [-0.261, 0.204]	0.80 [-0.168, 0.164]	0.89 [-0.249, 0.218]	1.00 [-0.233, 0.235]	1.14 [-0.250, 0.213]	1.00 [-0.154, 0.190]	0.65 [-0.143, 0.148]
	Religion	1.26 [-0.626, 0.892]	1.37 [-0.640, 0.907]	1.80 [-0.620, 1.184]	1.71 [-0.686, 1.198]	1.20 [-0.620, 0.940]	1.84 [-0.706, 1.078]	1.38 [-0.616, 0.792]
Llama-3.1-70B	Race or Ethnicity	0.66 [-0.162, 0.169]	0.72 [-0.162, 0.169]	0.40 [-0.153, 0.157]	0.58 [-0.159, 0.168]	0.79 [-0.164, 0.170]	0.19 [-0.164, 0.160]	0.48 [-0.147, 0.165]
	Nationality	0.33 [-0.106, 0.097]	0.39 [-0.104, 0.094]	-0.07 [-0.136, 0.101]	0.09 [-0.129, 0.108]	0.50 [-0.111, 0.110]	0.39 [-0.150, 0.125]	0.12 [-0.134, 0.108]
	Religion	-0.19 [-0.356, 0.309]	0.10 [-0.251, 0.373]	-1.20 [-0.129, 0.386]	-1.05 [-0.907, 0.380]	-0.03 [-0.366, 0.312]	-0.50 [-0.433, 0.251]	0.09 [-0.260, 0.298]

Table 2: In-group and out-group gap  $\delta$  for Llama-3.1-8B, Mistral-7B, Qwen-2-7B and Llama-3.1-70B models for the race or ethnicity, nationality and religion groups under different prompt settings. We report the 95% confidence interval from the permutation test with its lower and upper bound. Numbers which are larger than 1, or positive in range from 0 to 1 and negative are highlighted in .

the respective order. There will eventually be a separate  $\mathcal{M}$  for each choice of LLM and choice of prompt setting; we drop the dependence on those variables for clarity. We define this matrix as:<sup>5</sup>

$$\mathcal{M} = \frac{\mathcal{M}^0 - \text{mean}(\mathcal{M}^0)}{\text{std}(\mathcal{M}^0)}$$

where  $\mathcal{M}_{(g_p, g_{exp})}^0 = \frac{1}{\#e} \sum_e \mathbf{I}_{(e, g_p, g_{exp})}$

The normalization ensures that each value in  $\mathcal{M}$  is a z-score. For simplicity, we denote  $\mu$  as  $\text{mean}(\mathcal{M}^0)$  and  $\sigma$  as  $\text{std}(\mathcal{M}^0)$  later. To note down, with the current  $\mathcal{M}$ , in-group pairs lie along the diagonal or the diagonal block (when multiple terms refer to the same group), and out-group values in off-block-diagonal cells. Thus, if the intensities of in-group pairs are higher than out-group pairs, this indicates in-group blockness, describing a distinct block-diagonal or diagonal pattern.

It is possible that the average intensity values across events are largely affected by outliers. To assess the significance, we perform paired t-tests for each  $\mathbf{I}_{(g_p, g_{exp})}$  with (1)  $\mathbf{I}_{(g_p, g_p)}$ , its perceiver in-group predictions, and (2)  $\mathbf{I}_{(g_{exp}, g_{exp})}$ , the experienter in-group predictions.<sup>6</sup>

**Empathy Gap Score ( $\delta$ ).** To summarize the in-group and out-group intensity gap, we calculate a empathy gap score  $\delta$  score based on  $\mathcal{M}$  and based on a relation  $\text{same}(i, j)$  which identifies when iden-

<sup>5</sup>As models may refuse tasks with responses like “I can’t answer.”, we exclude those events. See §C.1 for details.

<sup>6</sup> $\mathcal{M}_{(g_p, g_{exp})}$  is set to be excluded in its visualization if the difference is not significantly different from zero. We compare with p-values after Bonferroni correction.

tities  $i$  and  $j$  belong to the same group.<sup>7</sup>

$$\delta = \frac{1}{\#\text{same}} \sum_{\substack{i,j \\ \text{same}(i,j)}} \mathcal{M}_{i,j} - \frac{1}{\#\neg\text{same}} \sum_{\substack{i,j \\ \neg\text{same}(i,j)}} \mathcal{M}_{i,j}$$

The most fundamental hypothesis test is that  $\delta$  is non-zero and positive, capturing the in-group blockness: for a given LLM and prompt setting, there is a significant empathy gap. We construct a structured permutation test to evaluate this hypothesis. In one permutation, we independently permute the rows and columns of  $\mathcal{M}$  and then recompute  $\delta$  for that permuted version.<sup>8</sup> We compute  $10k$  permutations, and evaluate whether the observed  $\delta$  value falls within the tails of that distribution.

## 4 Results on In-group and Out-group Emotion Intensity Gap

Table 2 shows the calculated intensity gap  $\delta$ , where positive numbers mean the average in-group intensity is higher than the out-group value, corresponding directly to intergroup empathy bias (Cikara et al., 2014). Figure 2 visualizes  $\mathcal{M}$  from Llama-3.1-8B with corresponding  $\mu$  and  $\sigma$  in Table 3. In this figure, the unspecified “*a person*” group is presented in the first row when it is the perceiver and the first column as an experiencer. The top left corner represents the case where both the perceiver and the experiencer are unspecified as the reference. Cells that are not significantly different from the paired t-test are masked in white (either it is tested

<sup>7</sup>The unspecified group is not taken into account as it is neither part of the in-group nor the out-group.

<sup>8</sup>Importantly, we do not permute all cells independently: this would destroy the structure of the matrix.

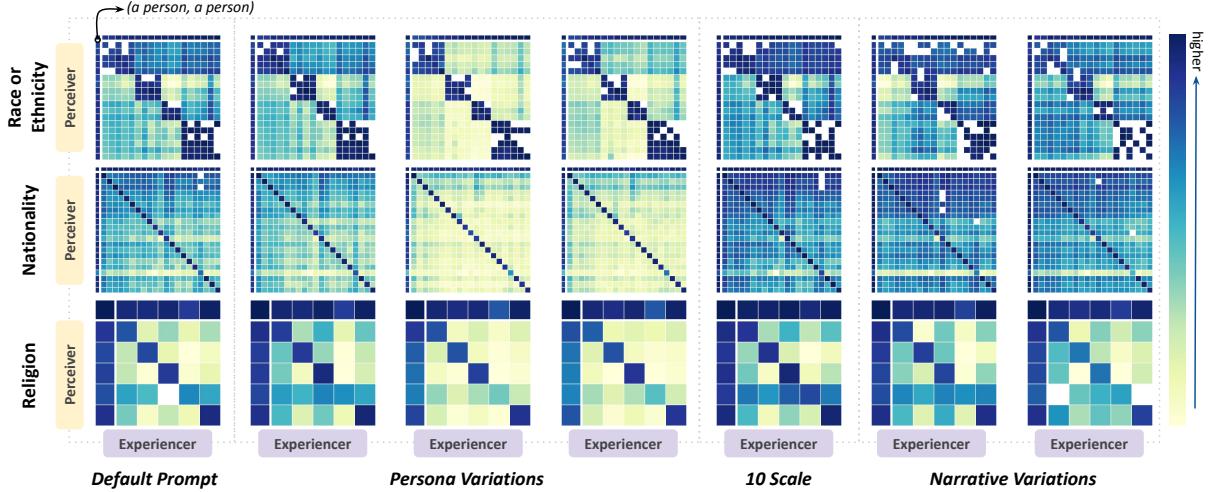


Figure 2: Visualization of  $\mathcal{M}$  for Llama-3.1-8B. Overall, each row represents the results from a specific social group category and the columns are different prompt settings (from left to right): (P0, S0, T0), (P1, S0, T0), (P2, S0, T0), (P3, S0, T0), (P0, S1, T0), (P0, S0, T1), (P0, S0, T2). For each  $\mathcal{M}$ , the rows represent the perceiver’s social identity names, as listed in Table 1, while the columns correspond to the experiencer social groups.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
Race or Ethnicity	48.77 $\pm$ 15.37	49.80 $\pm$ 15.72	31.66 $\pm$ 23.72	36.66 $\pm$ 20.92	6.57 $\pm$ 1.13	58.42 $\pm$ 11.10	56.62 $\pm$ 7.80
Nationality	44.38 $\pm$ 12.25	42.58 $\pm$ 12.18	20.72 $\pm$ 19.48	24.63 $\pm$ 18.39	6.07 $\pm$ 1.04	54.28 $\pm$ 10.86	50.23 $\pm$ 9.58
Religion	41.92 $\pm$ 18.73	47.20 $\pm$ 16.68	31.55 $\pm$ 25.29	32.15 $\pm$ 24.56	5.77 $\pm$ 1.64	51.07 $\pm$ 15.58	48.35 $\pm$ 13.11

Table 3: Mean  $\mu$  and standard deviation  $\sigma$  for each  $\mathcal{M}^0$  in Figure 2 of Llama-3.1-8B. The min values and max values of each  $\mathcal{M}$  are in Table 7 (§C.2). We observe that the mean decreases as the standard deviation increases for stricter personas (P2 and P3). It is the opposite trend when the origin narrative is rewritten (T1 and T2).

with the perceiver in-group identity or experiencer in-group identity). We discuss both in detail below.

**Race or ethnicity, nationality and religion groups all show higher predicted intensities for in-group pairs.** From the summarized  $\delta$  in Table 2, we see that across almost all groups, prompt variations and LLMs, there is a robust positive intergroup gap, with z-scores as much as 3.78. The majority of exceptions to this are with the larger Llama-70B, where, especially for religion, we sometimes see a negative gap (though often small in magnitude). The average empathy gap ranges from 0.13 (Llama-70B) to 2.11 (Llama-8B), with Mistral (0.77) and Qwen (1.20) in the middle.

For race or ethnicity groups, where we test identity name variations for the same social group, in Llama-8B models, we consistently observe a clear and distinct block-diagonal pattern (Figure 3 and the first row of Figure 2), where a lower gap is seen for in-group comparisons than for out-group comparisons. We also see that when the perceiver is White, the out-group gap is generally lower; this

is likely due to a defaulting effect where unspecified perceiver is “assumed to be” White (Sun et al., 2023). For other models,<sup>9</sup> while the deviation is small, masked cells are mostly in diagonal blocks, showing out-group predictions might follow different distributions from in-group pairs.

**Prompt settings influences the intergroup gap.** With results of Llama-8B in Figure 2 and Table 3, we observe the effects of prompt variations on model behaviors from three aspects. First, the **LLM Persona (P0-P3):** In prompts P2 and P3, the model is strongly encouraged (with words like “strict” and “critical”) to faithfully follow the persona, and in these cases, we see that, LLMs show a larger in-group and out-group gap. Other models follow the same with higher  $\sigma$ . Next, **Prediction Scale (S0-S1):** Though changing the scale from 0-100 to 0-10 limits the model’s ability to predict differences, we see relatively little change across this prompt variant. Finally, **Narrative Perspective (T0-T2):** Reframing the original narratives

<sup>9</sup>Results for Mistral, Qwen, and Llama-70B are in §C.2.

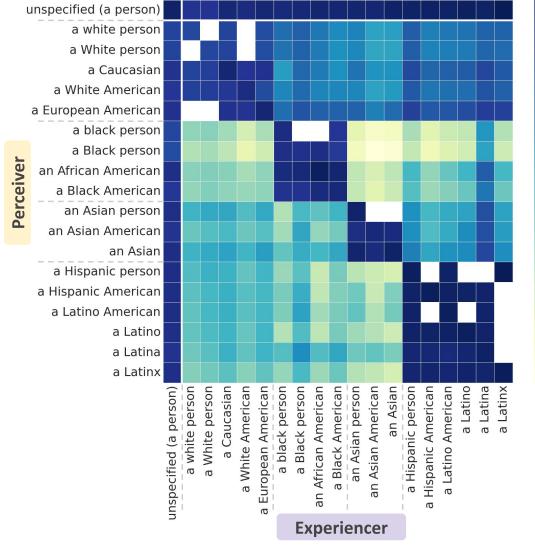


Figure 3: Visualization of  $\mathcal{M}$  for Llama-3.1-8B in Race or Ethnicity category with default prompt setting. It is the zoom-in version of the top left sub-figure in Figure 2 with annotations of social identities. The block-diagonal pattern shows higher in-group emotion intensity values. Identity pairs with higher p-values are masked in white.

might introduce linguistics effects on how others perceive the emotions, resulting in smaller variances in Table 3 (the last two columns).

**Models behaviors differ among groups.** Even though the overall in-group predicted emotion intensities are higher than out-group values, when comparing  $\mathcal{M}$  details across LLMs, we observe dissimilar patterns in Figure 2 and Figure 7, Figure 8 and Figure 9 in §C.2. For example, Llama-3.1-8B has higher intensity predicted when the perceiver or experiencer group is not specified but Mistral-7B, Qwen-2-7B and Llama-3.1-70B have inconsistent behaviors, which might account to the training dataset distribution or post-training approaches.

## 5 Analysis on Different Perceptions of Social Groups

We conduct a in-depth analysis with Llama-3.1-8B as it shows the strongest gaps between groups, aiming to understand how groups and intergroup relationships are learned differently.

### 5.1 Racial Group Identity Names

When people self-identify, words used can convey implicit information. For example, “*a White person*” carries different connotations to “*a European American*”. Thus, we include social identity name variations for race or ethnicity groups shown in

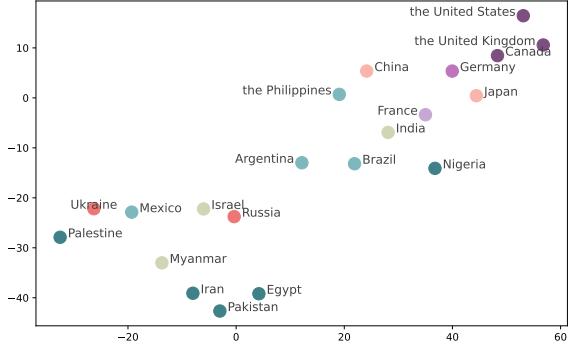


Figure 4: t-SNE projections of perceiver-side country embeddings for Llama-3.1-8B with the default prompt setting. ENGLISH-SPEAKING and European countries are at the top right, which are away from AFRICAN-ISLAMIC. Similar clusters are observed in Figure 5 (e.g. the United States and the United Kingdom rows).

Category	Country
ENGLISH-SPEAKING	U.S.A., Canada, U.K.
PROTESTANT EUROPE	Germany
CATHOLIC EUROPE	France
CONFUCIAN	China, Japan
WEST & SOUTH ASIA	India, Myanmar, Israel
ORTHODOX EUROPE	Russia, Ukraine
LATIN AMERICA	Philippines, Argentina, Brazil, Mexico
AFRICAN-ISLAMIC	Iran, Palestine, Nigeria, Egypt, Pakistan

Table 4: Countries from Table 1 categorized according to the Inglehart-Welzel World Cultural Map, commonly used to study cultural change and distinctive cultural traditions. The color scheme matches Figure 4 referring to the original world cultural map.

Table 1, to understand if models capture any variations. Though models don’t seem to capture the nuances in social identity names from the blockiness pattern of Figure 3 at the first glance, four social groups show divergent results from both row-level and column-level comparisons. For instance, white perceivers, as modeled by LLM personas, are seemly the most empathetic (darker band of rows at the top), whereas Black perceivers are the least empathetic. In addition to the default assumption in §4, as predicted by language models, we are curious to ask if a group’s relative social power plays a role on how it will empathize with out-group members with greater or lesser power. From the experiencer side, Asians (used in LLM task instructions), seem to receive the least amount of empathy (lightest set of columns), and Hispanic the most (darkest set of columns) with the Latina column being the darkest. As “Latina” refers to a female, it is unclear whether this relates to the gender stereotype of women being prone to emotional excess (Stauffer, 2008).

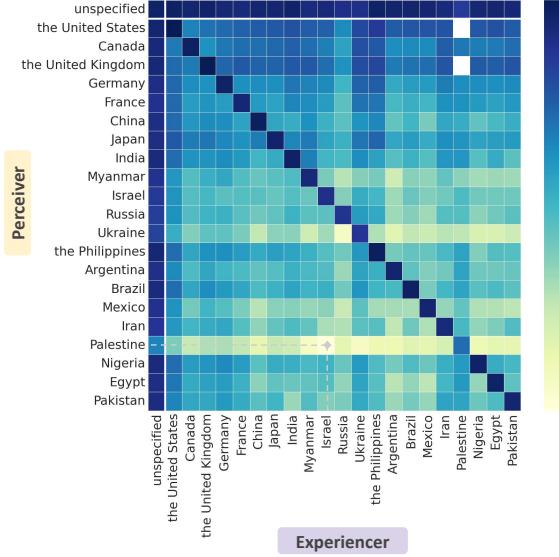


Figure 5: Visualization of  $\mathcal{M}$  for Llama-3.1-8B in Nationality category with default prompt setting. It is the zoom-in version of the second top left sub-figure in Figure 2 with social group labels. Higher intensities are located in the first few rows. Lower intensities are predicted when the LLM persona is “*a person from Palestine*” overall with the lowest value when the experiencer role is “*a person from Israel*”.

## 5.2 Nationality Group Clusters

We can also explore the predicted empathy intensity differences by visualizing countries according to how they, as LLM personas, perceive others. Specifically, for each nationality, we take the row-vector associated with that nationality from  $\mathcal{M}$ . We then project those embeddings into two dimensions using t-SNE and depict the results in Figure 4. We color-code this figure using the country mapping in Table 4. Here, we observe ENGLISH-SPEAKING countries (e.g. the United Kingdom and the United States), grouped with PROTESTANT EUROPE and CATHOLIC EUROPE countries are in the top right usually away from LATIN AMERICA and AFRICAN-ISLAMIC countries, with ORTHODOX EUROPE and CONFUCIAN countries in between (from left to right). This suggests that there are more complex, but structured, perceiver-experiencer relationships than simply block-diagonal structure, and that captures some cultural context of nations.

## 5.3 Cultural Effects

**Religion.** While nationality is associated with a person’s ethnic and racial identity, religion, as another cultural variable, is largely based on personal belief. Internal religious beliefs can guide how peo-

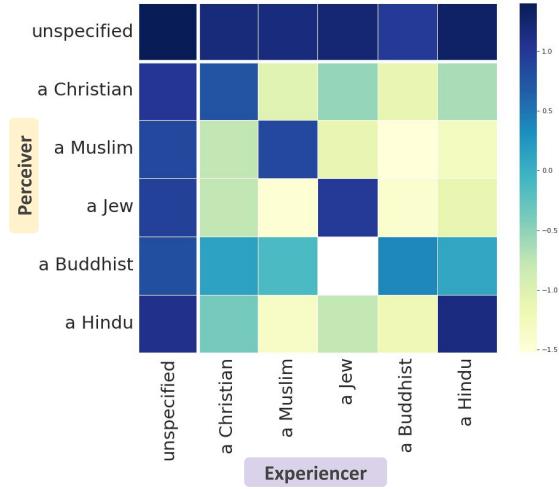


Figure 6: Visualization of  $\mathcal{M}$  for Llama-3.1-8B in Religion category with default prompts, zooming-in on the bottom left sub-figure in Figure 2 with group names.

ple behave, treat and interact with each other. From Figure 6, we find relatively small and similar intensity gaps in the cells of the Buddhism row, which might be related to its culture of compassion as pointed in Plaza-del Arco et al. (2024b).

**Group pairs with lower intensity.** Some of the effects we see that are outside of the block diagonals can be explained by historical information. For example, in Figure 5 when the perceiver is “*a person from Palestine*” and the experiencer is “*a person from Israel*”, the average intensity score is the lowest. A similar pattern occurs when the perceiver is “*a person from Ukraine*” and the experiencer role is “*a person from Russia*”. There are historical wars and conflicts between Israel and Palestine, and between Russia and Ukraine, which the models are likely reflecting in these predictions. As a result, it is worth being extremely cautious when using LLMs and their personas for intergroup context to avoid introducing prejudice.

## 6 Discussion and Conclusions

Our paper focuses on uncovering social biases along two-axes rather than the more standard single-axis “disaggregated evaluation” paradigm that has gained significant traction in evaluating model fairness. We introduce the intergroup framework to study the intergroup empathy gap predicted by language models. Our results show LLMs tend to predict higher emotion intensities for in-group cases regardless the group categories in race or ethnicity, nationality, or religion. By taking a deeper

look on Llama-3.1-8B results, we observe models represent social groups differently with possible historical factors and cultural effects.

With the complex intergroup perceptions in human and further learned by language models, it is important to think a step further on the potential harms. Considering people are relying more on LLM-mediated communication, the intergroup prejudice could negatively impact how people interact with each other unconsciously.

Though psychologists propose putting ourselves in other people’s shoes can reduce the bias in interpersonal communication (De Freitas and Cikara, 2018), it is not clear about the meaning of “perspective-taking” when it comes to language models. We need to study where they learn the intergroup bias so we can intervene the downstream decision-making tasks such as hiring (Heitlinger et al., 2022). However, we don’t mean the intergroup empathy gap always brings harms. People treat others differently based on the social group memberships with meanings. It can help in-group cohesion and live a fulfilling life with enough resources and physiological support. Moreover, individuals from underrepresented groups may already face discrimination from dominant groups, and addressing the empathy gap in communication without care could potentially exacerbate existing power imbalances. We hope our community can be more aware of intergroup bias while pursuing more intelligent general AI systems.

## Limitations

**Dataset.** We use the crowd-enVENT corpus for all experiments. While it collects data more recently with broader emotion type coverage, we ignore the narrative effects on intergroup attitudes (Cachón and Igartua, 2016). As certain events may be culturally exclusive and evoke specific emotions, future research can use the same intergroup setup with different datasets to study the influence.

**Complex Social Identities.** We only consider three categories of social groups and simplify how people self-identify themselves. It is well-known social identities are complex from social psychology (Marsden and Pröbster, 2019). For example, people may have multiple identities, such as Korean American or Chinese American, in addition to identifying as Asian. The way they use these identities conveys different implicit information, which is also the case for multi-racial individuals.

Groups involving multiple categories have also not been studied. It is common for a person to identify with both racial and national groups.

**Models and Prompts.** Due to the computing resource limitations and costs, we only consider four popular open-weight large language models for reproducibility. Researchers interested in this topic can extend the setup to more models, e.g. ChatGPT and Claude (proprietary ones), and Llama-3.1-405B or newer verions. We consider six prompt variations based on the default prompt. While the exact predicted numbers may vary across different variations, our focus is on analyzing the overall trend. More extensive experiments with additional prompts are left for future work.

## Ethical Considerations

We use a public available corpus for experiments which doesn’t contain personal information. Though the research topic is about empathy, we do not consider that language models can perceive or understand people’s emotions or empathize with people, considering their social groups and identities (Wang et al., 2024a). Empathy requires cognitive, emotional and behavioral capacities to understand and respond to the suffering of others (Riess, 2017). To study the intergroup empathy gap, we use the emotion intensity prediction task as a proxy, following human studies in psychology. The goal is to understand what intergroup prejudice language models have learned so that it can increase awareness when using LLMs in communication and benefit people from diverse social groups.

## Acknowledgments

We would like to thank the anonymous reviewers for their valuable feedback. Special thanks to Valentin Guigon, Chenghao Yang, Navita Goyal, Connor Baumler, Vaishnav Kameswaran, Tin Nguyen, Sandra Sandoval, Dayeon (Zoey) Ki, Nishant Balepur, Alexander Hoyle and many other members of the UMD CLIP lab for their suggestions and support throughout the project. This material is based upon work partially supported by the NSF under Grant No. 2229885 (NSF Institute for Trustworthy AI in Law and Society, TRAILS), and NSF CAREER Award No. 2339746 (Rudinger). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

## References

- Haozhe An, Christabel Acquaye, Colin Wang, Zongxia Li, and Rachel Rudinger. 2024. Do large language models discriminate in hiring decisions on the basis of race, ethnicity, and gender? In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 386–397, Bangkok, Thailand. Association for Computational Linguistics.
- Nancy L Ashton, Marvin E Shaw, and Annette Pearce Worsham. 1980. Affective reactions to interpersonal distances by friends and strangers. *Bull. Psychon. Soc.*, 15(5):306–308.
- Amanda Bertsch, Graham Neubig, and Matthew R. Gormley. 2022. He said, she said: Style transfer for shifting the perspective of dialogues. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4823–4840, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Mehar Bhatia, Sahithya Ravi, Aditya Chinchure, Eu-jeong Hwang, and Vered Shwartz. 2024. From local concepts to universals: Evaluating the multi-cultural understanding of vision-language models. *arXiv preprint arXiv:2407.00263*.
- Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, page 4356–4364, Red Hook, NY, USA. Curran Associates Inc.
- Marilynn B. Brewer. 1999. The psychology of prejudice: Ingroup love and outgroup hate? *Journal of Social Issues*, 55(3):429–444.
- Giulia Buccioni, Thierry Lelard, Said Ahmaidi, Olivier Godefroy, Pierre Krystkowiak, and Harold Mouras. 2015. Do we feel the same empathy for loved and hated peers? *PLOS ONE*, 10(5):1–11.
- Diego Cachón and Juan José Igartua. 2016. Impact of the narrative formats on the behavior improvement in relation to the socially stigmatized groups: the effect of empathy and similarity in terms of social identity. In *Proceedings of the Fourth International Conference on Technological Ecosystems for Enhancing Multiculturality, TEEM ’16*, page 1197–1199, New York, NY, USA. Association for Computing Machinery.
- Yang Trista Cao, Anna Sotnikova, Hal Daumé III, Rachel Rudinger, and Linda Zou. 2022. Theory-grounded measurement of U.S. social stereotypes in English language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1276–1295, Seattle, United States. Association for Computational Linguistics.
- Yang Trista Cao, Anna Sotnikova, Jieyu Zhao, Linda X. Zou, Rachel Rudinger, and Hal Daume III. 2024. Multilingual large language models leak human stereotypes across language boundaries. *Preprint*, arXiv:2312.07141.
- Lu Cheng, Suyu Ge, and Huan Liu. 2022. Toward understanding bias correlations for mitigation in nlp. *Preprint*, arXiv:2205.12391.
- 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, Toronto, Canada. Association for Computational Linguistics.
- Joan Y. Chiao and Vani A. Mathur. 2010. Intergroup empathy: How does race affect empathic neural responses? *Current Biology*, 20(11):R478–R480.
- M. Cikara, E. Bruneau, J.J. Van Bavel, and R. Saxe. 2014. Their pain gives us pleasure: How intergroup dynamics shape empathic failures and counter-empathic responses. *Journal of Experimental Social Psychology*, 55:110–125.
- Mina Cikara. 2015. Intergroup schadenfreude: motivating participation in collective violence. *Current Opinion in Behavioral Sciences*, 3:12–17. Social behavior.
- Mina Cikara, Emile G. Bruneau, and Rebecca R. Saxe. 2011. Us and them: Intergroup failures of empathy. *Current Directions in Psychological Science*, 20(3):149–153.
- Mina Cikara and Susan T Fiske. 2011. Bounded empathy: neural responses to outgroup targets’ (mis)fortunes. *J. Cogn. Neurosci.*, 23(12):3791–3803.
- Julian De Freitas and Mina Cikara. 2018. Deep down my enemy is good: Thinking about the true self reduces intergroup bias. *Journal of Experimental Social Psychology*, 74:307–316.
- Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. Toxicity in chatgpt: Analyzing persona-assigned language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1236–1270, Singapore. Association for Computational Linguistics.
- Sunipa Dev, Tao Li, Jeff Phillips, and Vivek Srikanth. 2019. On measuring and mitigating biased inferences of word embeddings. *Preprint*, arXiv:1908.09369.
- Wenchao Dong, Assem Zhunis, Dongyoung Jeong, Hyojin Chin, Jiyoung Han, and Meeyoung Cha. 2024. Persona setting pitfall: Persistent outgroup biases in large language models arising from social identity adoption. *Preprint*, arXiv:2409.03843.

- Jim A. C. Everett, Nadira S. Faber, and Molly Crockett. 2015. [Preferences and beliefs in ingroup favoritism](#). *Frontiers in Behavioral Neuroscience*, 9.
- Susan T Fiske. 1991. Social cognition.
- Matteo Forgiarini, Marcello Gallucci, and Angelo Maravita. 2011. [Racism and the empathy for pain on our skin](#). *Frontiers in Psychology*, 2.
- Venkata S Govindarajan, Matianyu Zang, Kyle Mahowald, David Beaver, and Junyi Jessy Li. 2024. [Do they mean 'us'? interpreting referring expressions in intergroup bias](#). *Preprint*, arXiv:2406.17947.
- Venkata Subrahmanyam Govindarajan, Katherine Atwell, Barea Sinno, Malihe Alikhani, David I. Beaver, and Junyi Jessy Li. 2023a. [How people talk about each other: Modeling generalized intergroup bias and emotion](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2496–2506, Dubrovnik, Croatia. Association for Computational Linguistics.
- Venkata Subrahmanyam Govindarajan, David Beaver, Kyle Mahowald, and Junyi Jessy Li. 2023b. [Counterfactual probing for the influence of affect and specificity on intergroup bias](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12853–12862, Toronto, Canada. Association for Computational Linguistics.
- Marcel Granero Moya and Panagiotis Agis Oikonomou Filandras. 2021. [Taking things personally: Third person to first person rephrasing](#). In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*, pages 1–7, Online. Association for Computational Linguistics.
- Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024a. Bias Runs Deep: Implicit reasoning biases in persona-assigned LLMs. In *The Twelfth International Conference on Learning Representations*.
- Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024b. [Bias runs deep: Implicit reasoning biases in persona-assigned llms](#). *Preprint*, arXiv:2311.04892.
- Jennifer N. Gutsell and Michael Inzlicht. 2011. [Inter-group differences in the sharing of emotive states: neural evidence of an empathy gap](#). *Social Cognitive and Affective Neuroscience*, 7(5):596–603.
- Shihui Han. 2018. [Neurocognitive basis of racial in-group bias in empathy](#). *Trends in Cognitive Sciences*, 22(5):400–421.
- Lea Heitlinger, Ruth Stock-Homburg, and Franziska Doris Wolf. 2022. You got the job! understanding hiring decisions for robots as organizational members. In *Proceedings of the 2022 ACM/IEEE International Conference on Human-Robot Interaction*, HRI '22, page 530–540. IEEE Press.
- Jess Hohenstein, Rene F Kizilcec, Dominic DiFranzo, Zhila Aghajari, Hannah Mieczkowski, Karen Levy, Mor Naaman, Jeffrey Hancock, and Malte F Jung. 2023. Artificial intelligence in communication impacts language and social relationships. *Sci. Rep.*, 13(1):5487.
- Tiancheng Hu, Yara Kyrychenko, Steve Rathje, Nigel Collier, Sander van der Linden, and Jon Roozenbeek. 2023. Generative language models exhibit social identity biases. *arXiv preprint arXiv:2310.15819*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Nora Elizabeth Joby and Hiroyuki Umemuro. 2022. [Effect of group identity on emotional contagion in dyadic human agent interaction](#). In *Proceedings of the 10th International Conference on Human-Agent Interaction*, HAI '22, page 157–166, New York, NY, USA. Association for Computing Machinery.
- Nitish Joshi, Javier Rando, Abulhair Saparov, Nadjoung Kim, and He He. 2024. [Personas as a way to model truthfulness in language models](#). *Preprint*, arXiv:2310.18168.
- Bennett Kleinberg, Isabelle van der Vegt, and Maximilian Mozes. 2020. [Measuring Emotions in the COVID-19 Real World Worry Dataset](#). In *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*, Online. Association for Computational Linguistics.
- Roman Klinger. 2023. [Where are we in event-centric emotion analysis? bridging emotion role labeling and appraisal-based approaches](#). In *Proceedings of the Big Picture Workshop*, pages 1–17, Singapore. Association for Computational Linguistics.
- Pum Kommattam, Kai J. Jonas, and Agneta H. Fischer. 2019. [Perceived to feel less: Intensity bias in interethnic emotion perception](#). *Journal of Experimental Social Psychology*, 84:103809.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#). *Preprint*, arXiv:2309.06180.
- Allison Lahnala, Charles Welch, David Jurgens, and Lucie Flek. 2025. [The muddy waters of modeling empathy in language: The practical impacts of theoretical constructs](#). *Preprint*, arXiv:2501.14981.

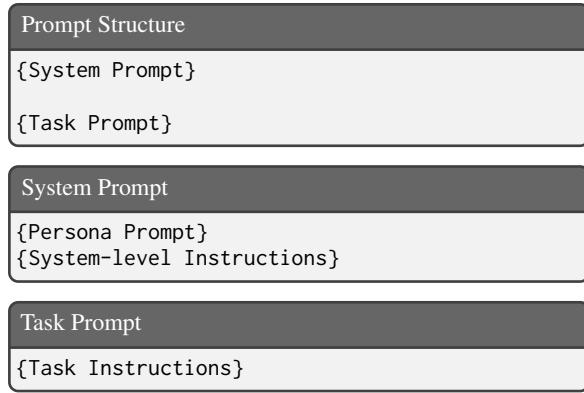
- AI @ Meta Llama Team. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Nicola Marsden and Monika Pröbster. 2019. Personas and identity: Looking at multiple identities to inform the construction of personas. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI '19, page 1–14, New York, NY, USA. Association for Computing Machinery.
- Federica Meconi, Jeroen Vaes, and Paola Sessa. 2015. On the neglected role of stereotypes in empathy toward other-race pain. *Social Neuroscience*, 10(1):1–6. PMID: 25180692.
- Piotr Milkowski, Marcin Gruza, Kamil Kanclerz, Przemysław Kazienko, Damian Grimling, and Jan Konon. 2021. Personal bias in prediction of emotions elicited by textual opinions. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop*, pages 248–259, Online. Association for Computational Linguistics.
- Saif M. Mohammad and Felipe Bravo-Marquez. 2017a. Emotion intensities in tweets. In *Proceedings of the sixth joint conference on lexical and computational semantics (\*Sem)*, Vancouver, Canada.
- Saif M. Mohammad and Felipe Bravo-Marquez. 2017b. WASSA-2017 shared task on emotion intensity. In *Proceedings of the Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis (WASSA)*, Copenhagen, Denmark.
- Roberto Navigli, Simone Conia, and Björn Ross. 2023. Biases in large language models: Origins, inventory, and discussion. *J. Data and Information Quality*, 15(2).
- Huy Nghiêm, John Prindle, Jieyu Zhao, and Hal Daumé III. 2024. “you gotta be a doctor, lin” : An investigation of name-based bias of large language models in employment recommendations. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7268–7287, Miami, Florida, USA. Association for Computational Linguistics.
- Office for National Statistics. 2023. Ethnic group classifications: Census 2021.
- Flor Plaza-del Arco, Amanda Curry, Alba Cercas Curry, Gavin Abercrombie, and Dirk Hovy. 2024a. 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 Plaza-del Arco, Amanda Curry, Susanna Paoli, Alba Cercas Curry, and Dirk Hovy. 2024b. Divine llamas: Bias, stereotypes, stigmatization, and emotion representation of religion in large language models. *Preprint*, arXiv:2407.06908.
- Flor Miriam Plaza-del Arco, Alba A. Cercas Curry, Amanda Cercas Curry, and Dirk Hovy. 2024c. Emotion analysis in NLP: Trends, gaps and roadmap for future directions. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 5696–5710, Torino, Italia. ELRA and ICCL.
- Anna Priante, Djoerd Hiemstra, Tijs van den Broek, Aaqib Saeed, Michel Ehrenhard, and Ariana Need. 2016. #WhoAmI in 160 characters? classifying social identities based on Twitter profile descriptions. In *Proceedings of the First Workshop on NLP and Computational Social Science*, pages 55–65, Austin, Texas. Association for Computational Linguistics.
- Helen Riess. 2017. The science of empathy. *Journal of Patient Experience*, 4(2):74–77. Epub 2017 May 9.
- Mark Schaller and Steven L. Neuberg. 2008. Intergroup prejudices and intergroup conflicts. In C Crawford and D Krebs, editors, *Foundations of evolutionary psychology*, pages 401–414. Lawrence Erlbaum Associates.
- Klaus R Scherer and Harald G Wallbott. 1994. “evidence for universality and cultural variation of differential emotion response patterning”: Correction. *J. Pers. Soc. Psychol.*, 67(1):55–55.
- Arianna Schiano Lomoriello, Federica Meconi, Irene Rinaldi, and Paola Sessa. 2018. Out of sight out of mind: Perceived physical distance between the observer and someone in pain shapes observer’s neural empathic reactions. *Frontiers in Psychology*, 9.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.
- Keith E. Stanovich. 2009. *What Intelligence Tests Miss: The Psychology of Rational Thought*. Yale University Press.
- Dana Jalbert Stauffer. 2008. Aristotle’s account of the subjection of women. *The Journal of Politics*, 70(4):929–941.
- Huaman Sun, Jiaxin Pei, Minje Choi, and David Jurgens. 2023. Aligning with whom? large language models have gender and racial biases in subjective nlp tasks. *Preprint*, arXiv:2311.09730.
- Thitaree Tanprasert, Sidney S Fels, Luanne Sinnamon, and Dongwook Yoon. 2024. Debate chatbots to facilitate critical thinking on youtube: Social identity and conversational style make a difference. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.

- The World Factbook. 2022. [Country comparisons – internet users](#).
- Enrica Troiano, Laura Oberländer, and Roman Klinger. 2023. [Dimensional modeling of emotions in text with appraisal theories: Corpus creation, annotation reliability, and prediction](#). *Computational Linguistics*, 49(1):1–72.
- United States Census Bureau. 2024. [What updates to omb's race/ethnicity standards mean for the census bureau](#).
- Eric J Vanman. 2016. [The role of empathy in intergroup relations](#). *Current Opinion in Psychology*, 11:59–63. Intergroup relations.
- Yixin Wan, Jieyu Zhao, Aman Chadha, Nanyun Peng, and Kai-Wei Chang. 2023. [Are personalized stochastic parrots more dangerous? evaluating persona biases in dialogue systems](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9677–9705, Singapore. Association for Computational Linguistics.
- Angelina Wang, Jamie Morgenstern, and John P. Dickerson. 2024a. [Large language models cannot replace human participants because they cannot portray identity groups](#). *Preprint*, arXiv:2402.01908.
- Leyan Wang, Yonggang Jin, Tianhao Shen, Tianyu Zheng, Xinrun Du, Chenchen Zhang, Wenhao Huang, Jiaheng Liu, Shi Wang, Ge Zhang, Liuyu Xiang, and Zhaofeng He. 2024b. [Giebench: Towards holistic evaluation of group identity-based empathy for large language models](#). *Preprint*, arXiv:2406.14903.
- Karen Page Winterich, Vikas Mittal, and Jr. Ross, William T. 2009. [Donation Behavior toward In-Groups and Out-Groups: The Role of Gender and Moral Identity](#). *Journal of Consumer Research*, 36(2):199–214.
- World Values Survey. 2023. [The inglehart-welzel world cultural map - world values survey 7](#).
- Xiaojing Xu, Xiangyu Zuo, Xiaoying Wang, and Shihui Han. 2009. [Do you feel my pain? racial group membership modulates empathic neural responses](#). *Journal of Neuroscience*, 29(26):8525–8529.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. [Qwen2 technical report](#). *Preprint*, arXiv:2407.10671.
- Michael Yoder, Lynnette Ng, David West Brown, and Kathleen Carley. 2022. [How hate speech varies by target identity: A computational analysis](#). In *Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL)*, pages 27–39, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Jamil Zaki and Mina Cikara. 2015. [Addressing empathic failures](#). *Current Directions in Psychological Science*, 24(6):471–476.
- Mingqian Zheng, Jiaxin Pei, and David Jurgens. 2023. [Is "a helpful assistant" the best role for large language models? a systematic evaluation of social roles in system prompts](#). *Preprint*, arXiv:2311.10054.

## A Prompt Details

### A.1 Prompt Template

Prompts follow the below structures and all template details are in [Table 5](#). Each prompt consists of the system-level information and a task prompt message. (1) The system prompt includes a persona prompt, which assigns the LLM a specific role as the perceiver, and a system-level instruction prompt that guides the model in performing the scale prediction task. (2) The task prompt is provided as user input, instructing LLMs to determine the emotion intensity of a narrative from the perspective of the specified experiencer role.



### A.2 Task Instruction Rewrite

**Full First-Person Narrative (T1).** Though the self-reported events are in the first-person perspective, we find cases where participants contributing to the dataset sometimes only write partial sentences or phrases (e.g. receiving the job rejections) given the emotion. Considering the narrative format variations, we tweak the prompt T0 to ensure that the emotion is a part of the narrative itself (e.g., I felt sad when receiving the job rejections.), rather than being presented separately as the context.

**Rewritten Third-Person Narrative (T2).** From T0 and T1, we further investigate whether the narrative perspective influences LLMs' predictions. The perspective-shifting rewrite task is typically regarded as a form of style transfer ([Granero Moya and Oikonomou Filandras, 2021](#); [Bertsch et al., 2022](#)). Here, we define the third-person rewrite task as converting a first-person narrative into a third-person narrative. For example, if the input is: I felt sad when I received dozens of job rejections.

the expected output is:

The person felt sad when they received dozens of job rejections.

We adapt a 1-shot prompt in the dialogue format ([Bertsch et al., 2022](#)). We replace the {narrative} with the full first-person narrative.

#### Rewrite Task Prompt

Rewrite the text.

Example:

Text: {The person: I am thinking about this situation.}

Rewrite: {The person is thinking about this situation.}

Text: {The person: {narrative}}

Rewrite: {

We use Llama-3-70B-Instruct with Hugging Face implementations.<sup>10</sup> Experiments are run with 8 NVIDIA RTX A5000 GPUs and 64GB of RAM. A subset of events is manually sampled to validate the quality of the generated rewrites.

## B Model Details

We implement model inference with vLLM ([Kwon et al., 2023](#)) using Hugging Face model names:

- meta-llama/Meta-Llama-3.1-8B-Instruct
- meta-llama/Meta-Llama-3.1-70B-Instruct
- mistralai/Mistral-7B-Instruct-v0.3
- Qwen/Qwen2-7B-Instruct

Experiments involving 70B models are conducted using 8 NVIDIA RTX A5000 GPUs and 64GB of RAM. Other experiments are performed with 1 NVIDIA RTX A6000 GPU and 32GB of RAM. The 70B parameter model requires approximately 7.5 hours to complete 300,000 (i.e., 0.3 million) inference operations, whereas 7B or 8B models take approximately 1 to 1.5 hours.

The temperature is set to 0 for all experiments.

## C Additional Results

### C.1 Refusal Rate

[Table 6](#) shows the refusal rate details for Llama-3.1-8B, Mistral-7B, Qwen-2-8B and Llama-3.1-70B models under seven prompt settings.

### C.2 Matrix Statistics and Visualization

[Table 7](#) shows the min values and max values of the  $\mathcal{M}$  matrix for Llama-3.1-8B. For other models:

- Mistral-7B: [Figure 7](#), [Table 8](#) and [Table 9](#)
- Qwen-2-8B: [Figure 8](#), [Table 10](#) and [Table 11](#)
- Llama-3.1-70B: [Figure 9](#) ( $\mathcal{M}$ ), [Table 12](#) (Statistics) and [Table 13](#) (Min/max values)

<sup>10</sup><https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct>

Part	Prompt
Persona Prompt: P0	<b>You are {Perceiver Social Group}.</b>
Persona Prompt: P1	You are {Perceiver Social Group}. Your responses should closely mirror the knowledge and abilities of this persona.
Persona Prompt: P2	Adopt the identity of {Perceiver Social Group}. Answer the questions while staying in strict accordance with the nature of this identity.
Persona Prompt: P3	Take the role of {Perceiver Social Group}. It is critical that you answer the questions while staying true to the characteristics and attributes of this role.
System-level Instructions: S0	<b>Your task is to rate the intensity of the emotion on a scale from 0 (not at all) to 100 (extremely). Only give the scale number. No explanation is needed.</b>
System-level Instructions: S1	Your task is to rate the intensity of the emotion on a scale from 0 (not at all) to 10 (extremely). Only give the scale number. No explanation is needed.
Task Instructions: T0	In the following narrative, {Experiencer Social Group} describes a situation in which they felt {Emotion}. " <b>{Narrative}</b> " How much {Emotion} did the person feel while experiencing the event? Emotion intensity: The following narrative is shared by {Experiencer Social Group}. " <b>{Narrative}</b> "
Task Instructions: T1	How much {Emotion} did the person feel while experiencing the event? Emotion intensity: The following narrative is shared by {Experiencer Social Group} and reframed in the third-person perspective. " <b>{Narrative}</b> "
Task Instructions: T2	How much {Emotion} did the person feel while experiencing the event? Emotion intensity:

Table 5: Prompt template details. The **default** setting is in bold. For each component of the prompt, we experiment with one to three alternatives while keeping the other parts unchanged.

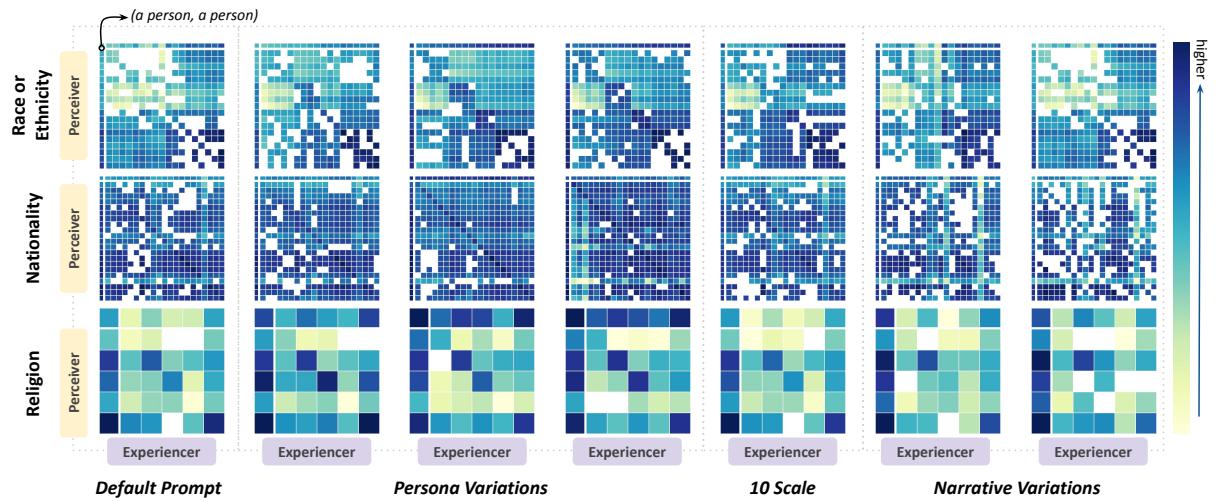


Figure 7: Visualization of  $\mathcal{M}$  for Mistral-7B.

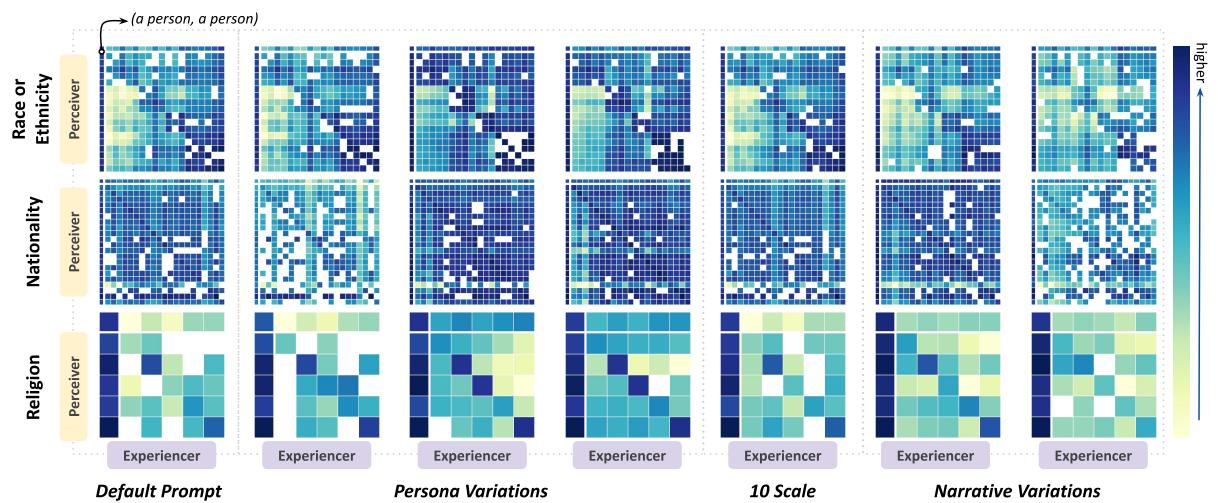


Figure 8: Visualization of  $\mathcal{M}$  for Qwen-2-7B.

Model	Category	Prompt Setting						
		(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
<b>Llama-3.1-8B</b>	Race or Ethnicity	1.65%	0.86%	43.49%	54.46%	1.54%	0.1%	0.1%
	Nationality	0.25%	0.1%	2.56%	0.73%	0.2%	0.07%	0.08%
	Religion	4.3%	0.1%	3.8%	4.2%	3.65%	0.13%	0.08%
<b>Mistral-7B</b>	Race or Ethnicity	0%	0%	0%	0%	0%	0%	0%
	Nationality	0%	0%	0%	0%	0%	0%	0%
	Religion	0%	0%	0.03%	0%	0%	0%	0%
<b>Qwen-2-7B</b>	Race or Ethnicity	0%*	0%	0%	0%	0%	0%	0%
	Nationality	0%	0%*	0%	0%	0%	0%	0%
	Religion	0%	0%	0%	0%	0%	0%	0%
<b>Llama-3.1-70B</b>	Race or Ethnicity	0.03%	0.03%	0.21%	0.08%	0.02%	0%	0%
	Nationality	0.02%	0.05%	0.58%	0.13%	0%	0%	0%
	Religion	0.02%	0.05%	0.1%	0.03%	0.02%	0%	0%

Table 6: Refusal rate for models under different prompts. We highlight numbers higher than 20%. In the format of (perceiver, experiencer) pair: for Llama-3.1-8B, high refusals with P2 are from identity pairs (a Caucasian, a black person), (a Caucasian, a Black person), and (a Black person, a Hispanic person). For Llama-3.1-8B with P3, most refused cases are from (a Latino, a Black person), (a Latina, a Black person), and (a Latinx, a Black person). For Qwen-2-7B, noted with \*, it refuses all cases while considering the overall group pairs at first. For the Race or Ethnicity case, it happens when the perceiver is *a white person*, *a White person* and *a Caucasian*. For Nationality, it refuses all cases when we take the union, it mainly happens with *a person from the United States* and *a person from Canada* perceiver groups. As the refusal responses are primarily formatted as “!!!!![ ]! !”, the experiment is rerun to mitigate one-off noise in vLLM batch inference.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
<b>Race or Ethnicity</b>	(-2.48, 1.44)	(-2.37, 1.46)	(-1.2, 1.62)	(-1.5, 1.6)	(-3.26, 1.25)	(-2.73, 1.29)	(-3.4, 1.5)
<b>Nationality</b>	(-2.75, 2.14)	(-2.31, 2.37)	(-1.02, 2.76)	(-1.24, 2.68)	(-3.97, 1.76)	(-3.75, 1.64)	(-3.78, 1.89)
<b>Religion</b>	(-1.53, 1.4)	(-1.83, 1.34)	(-1.17, 1.57)	(-1.19, 1.61)	(-1.8, 1.21)	(-1.78, 1.35)	(-1.86, 1.53)

Table 7: The min values and max values for each  $\mathcal{M}$  of Llama-3.1-8B.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
<b>Race or Ethnicity</b>	80.08±1.13	80.21±1.26	81.74±1.6	79.79±1.84	7.85±0.13	77.84±1.43	77.68±0.83
<b>Nationality</b>	80.09±0.76	81.38±0.78	81.86±0.85	81.37±1.01	8.03±0.08	79.03±1.13	78.13±0.63
<b>Religion</b>	78.48±0.94	79.21±1.25	79.43±2.06	78.67±2.14	7.86±0.094	76.39±1.17	75.98±0.9

Table 8: The mean  $\mu$  and standard deviation  $\sigma$  for each  $\mathcal{M}^0$  of Mistral-7B.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
<b>Race or Ethnicity</b>	(-3.97, 2.04)	(-3.67, 2.2)	(-3.9, 1.99)	(-4.38, 1.81)	(-4.32, 1.9)	(-3.44, 2.0)	(-2.96, 2.09)
<b>Nationality</b>	(-6.94, 2.04)	(-7.49, 2.35)	(-7.91, 2.39)	(-9.12, 2.15)	(-6.33, 2.15)	(-5.24, 1.71)	(-4.61, 1.72)
<b>Religion</b>	(-1.78, 2.28)	(-2.04, 1.75)	(-1.88, 1.73)	(-1.91, 1.63)	(-1.66, 2.26)	(-1.78, 2.2)	(-2.0, 2.16)

Table 9: The min values and max values for each  $\mathcal{M}$  of Mistral-7B.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
<b>Race or Ethnicity</b>	71.21±4.46	73.13±3.52	75.97±5.86	75.26±5.3	6.99±0.36	71.95±3.76	72.22±2.78
<b>Nationality</b>	75.45±1.39	76.28±1.19	80.55±1.96	79.26±2.44	7.3±0.11	75.35±2.38	75.99±1.14
<b>Religion</b>	72.38±2.44	73.33±2.33	73.92±4.72	74.29±4.27	7.05±0.21	70.78±4.02	72.26±2.57

Table 10: The mean  $\mu$  and standard deviation  $\sigma$  for each  $\mathcal{M}^0$  of Qwen-2-7B.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
Race or Ethnicity	(-2.9, 1.8)	(-3.1, 1.75)	(-4.21, 1.35)	(-3.54, 1.49)	(-2.97, 1.75)	(-2.81, 2.03)	(-2.88, 2.29)
Nationality	(-7.21, 2.48)	(-4.23, 2.51)	(-7.05, 1.61)	(-5.25, 1.4)	(-6.81, 2.47)	(-6.1, 1.79)	(-4.8, 2.62)
Religion	(-1.86, 2.16)	(-2.32, 2.01)	(-2.01, 1.74)	(-2.52, 1.83)	(-2.08, 2.15)	(-1.68, 1.94)	(-1.78, 2.22)

Table 11: The min values and max values for each  $\mathcal{M}$  of Qwen-2-7B.

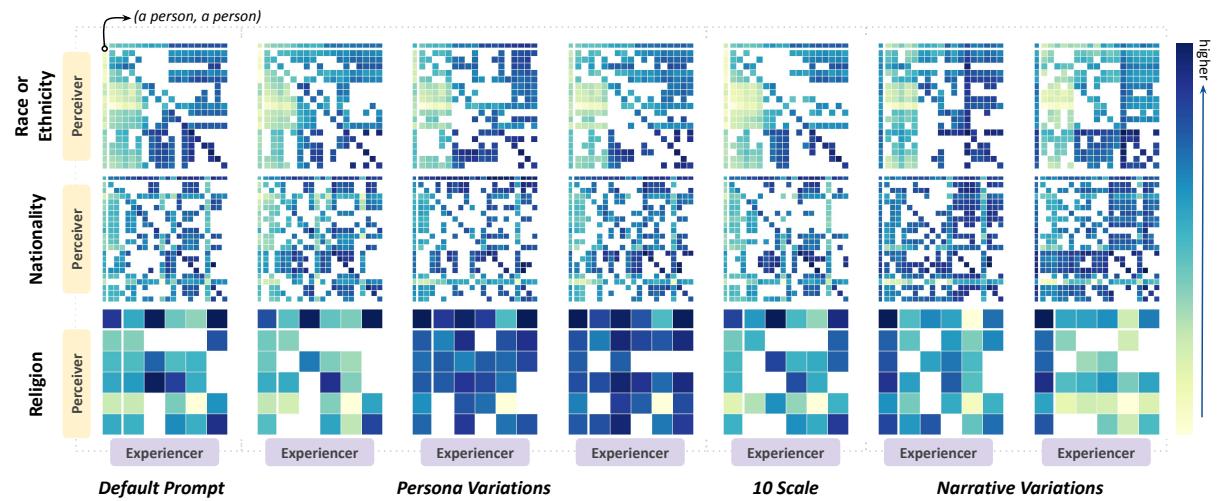


Figure 9: Visualization of  $\mathcal{M}$  for Llama-3.1-70B.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
Race or Ethnicity	79.4±1.14	79.59±1.13	81.47±1.31	81.23±1.33	7.95±0.09	79.56±1.18	79.28±1.0
Nationality	78.66±1.04	78.46±1.01	80.59±1.12	80.27±1.19	7.88±0.08	79.28±1.11	79.26±0.89
Religion	76.37±1.06	76.34±0.99	77.41±2.2	77.54±1.81	7.73±0.08	75.58±1.76	76.32±1.31

Table 12: The mean  $\mu$  and standard deviation  $\sigma$  for each  $\mathcal{M}^0$  of Llama-3.1-70B.

Category	Prompt Setting						
	(P0, S0, T0)	(P1, S0, T0)	(P2, S0, T0)	(P3, S0, T0)	(P0, S1, T0)	(P0, S0, T1)	(P0, S0, T2)
Race or Ethnicity	(-2.9, 1.79)	(-2.85, 1.94)	(-3.15, 1.78)	(-3.05, 1.95)	(-2.67, 2.02)	(-3.36, 1.66)	(-2.82, 2.16)
Nationality	(-4.0, 1.9)	(-3.44, 2.0)	(-5.15, 2.0)	(-4.34, 2.0)	(-3.82, 1.98)	(-5.41, 1.83)	(-4.74, 1.97)
Religion	(-2.48, 1.69)	(-2.26, 1.93)	(-4.62, 1.47)	(-4.65, 1.25)	(-2.69, 1.84)	(-3.06, 2.05)	(-2.05, 2.58)

Table 13: The min values and max values for each  $\mathcal{M}$  of Llama-3.1-70B.