# Large Language Models Portray Socially Subordinate Groups as More Homogeneous, Consistent with a Bias Observed in Humans

Messi H.J. Lee
hojunlee@wustl.edu
Division of Computational and Data
Sciences, Washington University in St.
Louis
St. Louis, Missouri, USA

Jacob M. Montgomery
jacob.montgomery@wustl.edu
Department of Political Science,
Washington University in St. Louis
St. Louis, Missouri, USA

Calvin K. Lai
calvinlai@wustl.edu
Department of Psychological & Brain
Sciences, Washington University in St.
Louis
St. Louis, Missouri, USA

## ABSTRACT

Large language models (LLMs) are becoming pervasive in everyday life, yet their propensity to reproduce biases inherited from training data remains a pressing concern. Prior investigations into bias in LLMs have focused on the association of social groups with stereotypical attributes. However, this is only one form of human bias such systems may reproduce. We investigate a new form of bias in LLMs that resembles a social psychological phenomenon where socially subordinate groups are perceived as more homogeneous than socially dominant groups. We had ChatGPT, a state-of-the-art LLM, generate texts about intersectional group identities and compared those texts on measures of homogeneity. We consistently found that ChatGPT portrayed African, Asian, and Hispanic Americans as more homogeneous than White Americans, indicating that the model described racial minority groups with a narrower range of human experience. ChatGPT also portrayed women as more homogeneous than men, but these differences were small. Finally, we found that the effect of gender differed across racial/ethnic groups such that the effect of gender was consistent within African and Hispanic Americans but not within Asian and White Americans. We argue that the tendency of LLMs to describe groups as less diverse risks perpetuating stereotypes and discriminatory behavior.

## CCS CONCEPTS

• **Applied computing** → **Psychology**; • **Computing methodologies** → **Natural language processing**.

## KEYWORDS

Large Language Models, AI Bias, Homogeneity Bias, Perceived Variability, Stereotyping

## 1 INTRODUCTION

In recent years, the examination of bias in Artificial Intelligence (AI) has garnered significant attention, with multiple studies spotlighting biases in AI systems designed for real-world decision-making [e.g., 10, 18, 19]. For instance, Buolamwini and Gebru [10] showed that commercial gender classification systems, used in various sectors like marketing, entertainment, security, and healthcare, achieved higher accuracy for lighter-skinned individuals than darker-skinned individuals, and that the disparity was most pronounced within darker-skinned females with error rates high as 34.7% (as opposed to 0.3% of lighter-skinned males). This study, along with many others, demonstrated that AI systems, contrary to the expectation that they would be impartial and immune to biases, could show performance disparities for specific groups and reproduce, or even amplify, human biases.

Natural language processing (NLP) systems are similarly vulnerable to bias. Since the seminal works of Bolukbasi et al. [6] and Caliskan et al. [11] documenting human-like biases within word embedding models, a wide array of studies have found biases within models for coreference resolution [49], text classification [15], machine translation [38, 46], and text generation [1, 34], among many others. For example, Lucy and Bamman [34] showed that GPT-3 would write stories related to family, emotions, and body parts when asked to write about a feminine character whereas it would write stories related to politics, war, sports, and crime when asked to write about a masculine character. Another work by Abid et al. [1] showed that GPT-3 would associate Muslims with violence when performing text completions. These studies highlighted the role Large Language Models (LLMs) could play in reproducing and amplifying stereotypical trait associations in their generated content.

### 1.1 Biases beyond trait association

The above studies not only underscore the potential for LLMs to reproduce and amplify stereotypical trait associations, but they also prompt researchers to question whether LLMs reproduce other human-like biases. One type of bias that remains unexplored in LLMs is perceived homogeneity of groups - the tendency to perceive some social groups as less diverse/more homogeneous compared to others. This bias was first studied within the context of intergroup relations where social psychologists found that people tend to perceive members of their outgroup as more homogeneous than members of their ingroup [30]. Subsequently, the phenomenon was documented across a wide variety of social distinctions including gender [36], age [29], race/ethnicity [2], college majors [37], and political orientation [39]. However, further exploration

revealed that differences in the perceived homogeneity of ingroups and outgroups may instead be attributable to the relative social status and power of groups [22–24, 32, 33]. These studies found that members of socially dominant groups perceived their outgroup(s) as more homogeneous than the ingroup (in line with the typical outgroup homogeneity effect), but that members of socially subordinate groups would perceive their ingroup(s) as more homogeneous than the socially dominant outgroup. Together, these effects suggest that humans have a general tendency to perceive socially subordinate groups as more homogeneous than socially dominant groups.

Perceived homogeneity (or variability) of groups is a form of stereotyping that has strong implications for prejudice and discrimination. Studies show that viewing a group as more variable reduces other forms of stereotyping [25, 43], prejudice, and discrimination [7, 20]. As LLMs become increasingly involved in everyday life, it is essential to understand if they perpetuate biases related to perceived homogeneity as they may influence users' perceptions and attitudes towards groups. This investigation is part of a broader discussion on erasure within Natural Language Processing [NLP; 16, 17], which highlights the lack of adequate representation of social groups in NLP systems. Homogeneous representations of subordinate groups in LLM outputs, or *homogeneity bias*, not only undermine the rich and diverse identities of these groups but also reinforce existing social hierarchies.

## 1.2 Homogeneous narratives of marginalized groups in LLMs

Recent works in the LLM literature, such as Cheng et al. [12] and Cheng et al. [13], have highlighted LLMs' tendencies to essentialize and produce positive yet homogeneous narratives of marginalized groups in personas, written descriptions of an individual who identifies with a given social group identity (e.g., "Imagine you are an Asian woman. Describe yourself."). Cheng et al. [13] measure the extent to which these descriptions focus on groups' defining characteristics, often linked to stereotypes, in a manner akin to "stereotype endorsement," one of three types of measures used to study the outgroup homogeneity effect [35]. Building on this, we introduce a new method to assess homogeneity in group representations, akin to "perceived similarity," which quantifies the degree of similarity in these representations. Furthermore, we extend our analysis to text formats more aligned with everyday use of LLMs (e.g., stories), underscoring the pervasive harm of homogeneity bias. Our findings indicate that homogeneity bias affects not only the content but also the manner in which the narratives are conveyed.

## 1.3 This work

In this work, we empirically test whether LLMs exhibit bias akin to human perceptions of group homogeneity through an experiment using ChatGPT. We had ChatGPT generate texts about eight different intersectional groups. We looked at four racial/ethnic groups - African, Asian, Hispanic, and White Americans - where White Americans were identified as the dominant racial/ethnic group [51], and we looked at two gender groups - men and women - where men were identified as the dominant gender group [47]. If LLMs reproduce this human-like bias, we would expect LLMs to describe

members of the socially subordinate group as more homogeneous than those of the socially dominant group. We formalize our pre-registered research questions[1] as follows:

**Research Question 1.** *Does ChatGPT depict U.S. racial/ethnic minority groups (African, Asian, and Hispanic Americans) as more homogeneous compared to the U.S. racial/ethnic majority group (White Americans)?*

**Research Question 2.** *Does ChatGPT depict the gender minority group (women) as more homogeneous compared to the gender majority group (men)?*

**Research Question 3.** *Is the effect of gender on the homogeneity of text generated by ChatGPT consistent across racial/ethnic groups?*

## 2 METHOD

### 2.1 Data

We created a collection of writing prompts asking ChatGPT to write texts about eight intersectional group identities. We included four racial/ethnic groups - African, Asian, Hispanic, and White Americans - and two gender groups - men and women. To generate a wide range of comparable content, we considered a variety of text formats such as stories, character descriptions, and biographies. To control for text length, we limited generated text to 30 words.[2] The prompts read, "Write a 30-word [ story about / character description of / biography of / introduction of / social media profile of / synopsis for / narrative of / self-introduction of / tragic story about / funny story about / romantic story about / horror story about / dramatic story about ] a(n) [ African / Asian / Hispanic / White ] American [ man / woman ]."

We used the OpenAI API, specifically employing the `gpt-3.5 -turbo` model (as of 25 July 2023) to obtain 500 text completions for each prompt. The decision to collect 500 completions stemmed from pilot tests suggesting that a smaller number of completions (i.e., 10 or 100) lead to more instability in our estimates. We used the default parameters of the API,[3] but made two exceptions: the `n` parameter, which determines the number of text completions per API request, and the `role` of the system that determines the model's behavior (set to "chatbot").[4] To ensure data quality, we did a keyword-based query to identify and remove 50 out of 52,000 instances where ChatGPT refused to generate the requested texts.[5]

### 2.2 Measure of text homogeneity

We assessed text homogeneity by calculating the pairwise cosine similarity between sentence embeddings of texts generated for each group. These embeddings are numeric vectors in a multidimensional space that encode the semantic and syntactic information of sentences [14]. We obtained these embeddings using the second-to-last layer of the BERT-base-uncased model, referred to below

---

[1] https://osf.io/kxz6b/
[2] ChatGPT did not strictly follow the length requirement. The texts had an average length of 26.61 words (*SD* = 2.70).
[3] https://platform.openai.com/docs/guides/gpt-chat-completions-api
[4] We gathered data in four separate batches, with `n`s set to 128, 128, 128, and 116 as the API could only process up to 128 generations in each request.
[5] We provide a breakdown of non-compliant completions by race/ethnicity, gender, and text format in Section A.4 of the Supplementary Materials. These non-compliant completions were replaced with new ones.

as BERT$_{-2}$, following our pre-registered analysis plan. This choice aligned with the default configuration of the `text` R package [R Version 4.3.1; 26] and reflected the fact that upper layers (i.e. close to last) of the embedding model tend to provide more contextualized representations of language [21].

We conducted four sets of additional analyses to evaluate the robustness of our findings to alternative approaches for measuring similarity (these were not pre-registered). We used (1) the third-to-last layer of BERT (BERT$_{-3}$), (2) the second-to-last layer of the larger RoBERTa-base model [31, RoBERTa$_{-2}$], (3) the third-to-last layer of RoBERTa (RoBERTa$_{-3}$), and (4) three pre-trained Sentence-BERT models with highest average performance on sentence encoding tasks [40]: `all-mpnet-base-v2`, `all-distilroberta-v1`, and `all-MiniLM-L12-v2`.

After encoding the ChatGPT-generated texts into sentence embeddings, we calculated the cosine similarity between all pairs of the sentence embeddings that were induced for each of the prompts. Cosine similarity is calculated by taking the dot product of two sentence embeddings and dividing it by the product of their magnitudes. The value can range from -1 to 1, where 1 indicates that the two sentences are perfectly identical and where -1 indicates that the two sentences are completely dissimilar. We then standardized this measure for interpretability (subtracting the mean and dividing by the standard deviation). Table 1 shows the *most similar* and *least similar* pairs of texts according to the standardized cosine similarity values computed using BERT$_{-2}$. These examples provide some face validity to our measurement strategy as the first sentence pair largely conveys the same message while the second pair does not. To see if this generalizes, we present ten random sentence pairs in Table A1 of the Supplementary Materials. These examples again provide strong face validity for our measurement strategy, with high-scoring pairs appearing to be far more similar than low-scoring pairs. As we generated 500 texts for each prompt, there were 124,750 pairs of sentence embeddings, and hence 124,750 cosine similarity measurements corresponding to each prompt.

### 2.3 Testing group differences

Following the pre-registered analysis plan, we used linear mixed-effects models with functions from the `lme4` [3] and `lmerTest` [27] R packages. In the models, we included race/ethnicity, gender, and their interactions as fixed effects and text format as random intercepts. Text format was included as random intercepts instead of random slopes because we expected the cosine similarity baseline to vary across text formats,[6] but we did not expect the magnitude and direction of race/ethnicity and gender to vary across text format.[7]

We also fitted additional un-pre-registered models to facilitate interpretation of race/ethnicity and gender fixed effects in the presence of interactions [8]. We fitted mixed-effects models where (1) race/ethnicity was the only fixed effect ("Race/Ethnicity model"), (2) gender was the only fixed effect ("Gender model"), and (3) race/ethnicity and gender were both fixed effects ("Race/Ethnicity &

Gender model"). These models allowed for easier interpretation and led to the same substantive conclusions. Subsequently, we used the pre-registered mixed-effects model ("Interaction model") to interpret the interaction effect.

We used the `afex` R package [45] to conduct likelihood-ratio tests to determine if the models including the fixed effects of race/ethnicity, gender, and their interactions provided better fits for the data than those without. To determine the magnitude and direction of race/ethnicity and gender, we examined the summary outputs of the Race/Ethnicity and Gender models. Finally, to examine the interaction effects, we used the `emmeans` R package [28] to conduct pairwise comparisons of estimated marginal means between gender groups within the same racial/ethnic groups. In all models, White Americans and men served as reference categories.[8]

## 3 RESULTS

In Table 2, we present the means and standard deviations of the standardized cosine similarity values for the eight intersectional groups, computed using BERT$_{-2}$.

### 3.1 Main effect of race/ethnicity

ChatGPT-generated texts about the subordinate racial/ethnic groups were more homogeneous than those about the dominant racial/ethnic group (see Figure 1). The Race/Ethnicity model (Column 1 in Table 3) showed that the standardized cosine similarity values of African, Asian, and Hispanic Americans were each 0.33 ($SE < .001$, $t(12,973,984) = 508.81$), 0.31 ($SE < .001$, $t(12,973,984) = 478.74$), and 0.18 ($SE < .001$, $t(12,973,984) = 275.05$) standard deviations greater than those of White Americans. In addition, the likelihood-ratio test showed that the model including race/ethnicity provided a better fit for the data than that without it, as indicated by the chi-squared statistics for the analysis using BERT$_{-2}$ ($\chi^2(3) = 326701.07$, $p < .001$; see Table A3). These findings replicated across all six alternative measurement strategies. For results of the likelihood ratio tests, see Table A3, and for summary outputs of the mixed effects models, see Tables A5-A10.



**Figure 1: Mean standardized cosine similarity values of the four racial/ethnic groups using BERT$_{-2}$. Error bars were omitted as confidence intervals were all smaller than 0.001.**

---

[6]Text formats like self-introduction, for example, may be more similar to each other than other text formats given that self-introductions are likely to share a common structure or content that constitutes an introduction.

[7]When fitting linear mixed-effects models, we turned off derivative calculations that could slow down the model fitting process and used the `nmkbw` optimizer made available by the `lme4` R package.

[8]Code is available at https://github.com/lee-messi/Homogeneity-Bias-in-LLMs

**Table 1: Pairs of sentences with the highest and lowest standardized cosine similarity values among stories written about African American men. The cosine similarity values were calculated using BERT$_{-2}$.**

| Sentence 1 | Sentence 2 | Std. Cos. Sim. |
|---|---|---|
| In a world divided by prejudice, he shattered stereotypes with his compassionate heart, empowering others to rise above discrimination and embrace unity. | In a world divided by prejudice, he defied stereotypes with his intelligence and compassion, inspiring others to rise above ignorance and embrace unity. | 1.57 |
| He closed his eyes and took a deep breath, feeling the weight of history on his shoulders. With determination, he stepped forward, ready to redefine his legacy. | An African American man woke up to a world where color no longer mattered, and everyone saw the brilliance in every shade of skin. | −4.98 |

**Table 2: Descriptive statistics of the standardized cosine similarity values for the eight intersectional groups. Cosine similarity computations were performed using BERT$_{-2}$ and were then standardized for better interpretability.**

| Race/Ethnicity | Gender | N | Mean | St. Dev. |
|---|---|---|---|---|
| African Americans | Men | 124,750 | 0.12 | 0.79 |
| | Women | 124,750 | 0.13 | 0.86 |
| Asian Americans | Men | 124,750 | 0.10 | 0.83 |
| | Women | 124,750 | 0.11 | 0.87 |
| Hispanic Americans | Men | 124,750 | -0.09 | 1.34 |
| | Women | 124,750 | 0.04 | 1.25 |
| White Americans | Men | 124,750 | -0.21 | 0.89 |
| | Women | 124,750 | -0.21 | 0.95 |



**Figure 2: Standardized cosine similarity values of the two gender groups using BERT$_{-2}$. Error bars were omitted as confidence intervals were all smaller than 0.001.**

## 3.2 Main effect of gender

ChatGPT-generated texts about the subordinate gender group (i.e., women) were also more homogeneous than those about the dominant gender group (men), although the differences were modest (see Figure 2). The Gender model in Table 3 showed that the cosine similarity values of women were 0.037 ($SE < .001$, $t(12,973,986) = 78.68$) standard deviations greater than those of men.[9] Furthermore, the likelihood-ratio test found that the model including the gender term provided a better fit for the data than that without it, as indicated by the chi-squared statistics for the analysis using BERT$_{-2}$ ($\chi^2(1) = 6352.47$, $p < .001$; see Table A3). These findings replicated across all six alternative measurement strategies. For results of the likelihood ratio tests, see Table A3, and for summar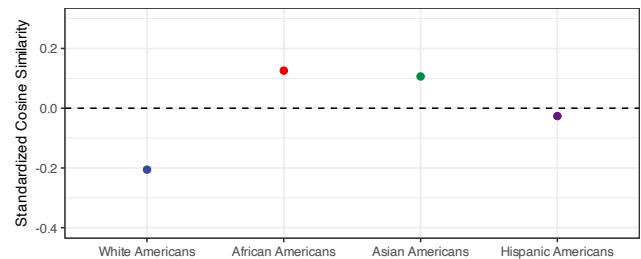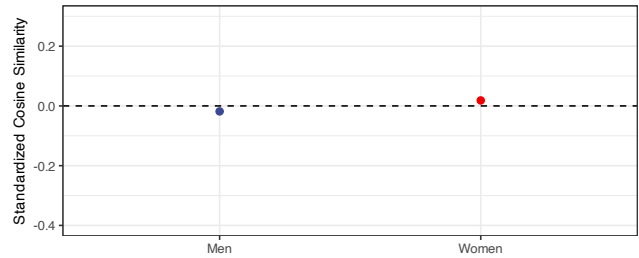y outputs of the mixed effects models, see Tables A5-A10. However, we note that, although statistically significant, these results indicated that the impact of gender was substantially smaller than that of race/ethnicity.
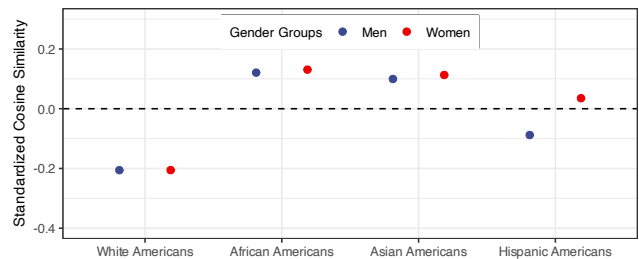
## 3.3 Interaction effect

The effect of gender on the homogeneity of ChatGPT-generated text differed between racial/ethnic groups. Pairwise comparisons

of estimated marginal means revealed that African, Asian, and Hispanic American women each held greater cosine similarity values than their male counterparts ($z$s = 10.79, 14.54, 133.86, $p$s < .001), but there was no significant difference between White American men and women ($z = 0.23$, $p = .82$; see Table A4 and Figure 3). The likelihood-ratio test found that the model including the interaction term provided a better fit for the data than that without it, as indicated by the chi-squared statistics for the analysis using BERT$_{-2}$ ($\chi^2(3) = 11888.15$, $p < .001$; see Table A3).



**Figure 3: Standardized cosine similarity values of all eight intersectional groups using BERT$_{-2}$. Error bars were omitted as confidence intervals were all smaller than 0.001.**

We observed slight variations in the effects of gender within individual racial/ethnic groups when alternative measurement strategies involving BERT and RoBERTa were used (see Figure 4). Examining the results in Table A4, African American women held greater cosine similarity values than their male counterpart ($z$s = 15.34,

---

[9]The base term for gender in the Interaction model (Column 4 of Table 3) was not significant, but this does not mean that gender had no effect. Rather, this indicates that gender had no measurable effect within White Americans (the reference category). We discuss this further in the next section.

**Table 3: Summary output of mixed effects models using cosine similarity values from BERT$_{-2}$. Positive coefficients indicate greater pairwise cosine similarity and thus more homogeneity compared to the baseline categories - White Americans and men.**

| | BERT$_{-2}$ | | | |
|---|---|---|---|---|
| | Race/Ethnicity model | Gender model | Race/Ethnicity, Gender model | Interaction model |
| **Fixed Effects** | | | | |
| Intercept | −0.21 (0.16) | −0.018 (0.16) | −0.22 (0.16) | −0.21 (0.16) |
| African Americans | 0.33* (0.00065) | | 0.33* (0.00065) | 0.33* (0.00092) |
| Asian Americans | 0.31* (0.00065) | | 0.31* (0.00065) | 0.31* (0.00092) |
| Hispanic Americans | 0.18* (0.00065) | | 0.18* (0.00065) | 0.12* (0.00092) |
| Women | | 0.037* (0.00047) | 0.037* (0.00046) | 0.00021 (0.00092) |
| African Americans × Women | | | | 0.0097* (0.0013) |
| Asian Americans × Women | | | | 0.013* (0.0013) |
| Hispanic Americans × Women | | | | 0.12* (0.0013) |
| **Random Effects ($\sigma^2$)** | | | | |
| Text Format Intercept | 0.32 | 0.32 | 0.32 | 0.32 |
| Residual | 0.69 | 0.71 | 0.69 | 0.69 |
| Observations | 12,974,000 | 12,974,000 | 12,974,000 | 12,974,000 |
| Log likelihood | −15,985,323 | −16,145,340 | −15,982,157 | −15,976,230 |

$^*p < .001$

82.55, 44.27, $p$s $< .001$), Asian American women held greater cosine similarity values than their male counterpart ($z$s = 34.32, 100.39, 72.79, $p$s $< .001$), and Hispanic American women held greater cosine similarity values than their male counterpart ($z$s = 142.07, 141.82, 145.79, $p$s $< .001$). However, unlike the pre-registered analysis reported in Table A4, White American women also held greater cosine similarity values than their male counterpart ($z$s = 22.61, 117.75, 99.70, $p$s $< .001$).[10]

We observed more variations in the effects of gender within individual racial/ethnic groups when alternative measurement strategies involving Sentence-BERT were used. Consistent with the pre-registered analysis, African American women held greater cosine similarity values than their male counterpart ($z$s = 98.34, 95.25, 64.65, $p$s $< .001$), and Hispanic American women held greater cosine similarity values than their male counterpart ($z$s = 352.72, 351.10, 224.90,

$p$s $< .001$). However, the direction of the effect of gender within Asian Americans differed across models ($z$s = 5.81, −40.29, −47.15, $p$s $< .001$). Similarly, the direction of the effect of gender within White Americans differed across models ($z$s = 4.61, −45.44, −52.52, $p$s $< .001$). All in all, the effect of gender was consistent in one direction within African and Hispanic Americans but not within Asian and White Americans.[11]

### 3.4 Homogeneity bias and topical alignment

In Section A.2 of the Supplementary Materials, we conducted two un-pre-registered follow-up studies and an exploratory analysis to unpack the source of homogeneity bias as measured from cosine similarity of sentence embeddings. We explored whether topical alignment, defined as the frequency of shared topics in texts about specific groups, might account for the observed homogeneity bias.

---

[10]The likelihood-ratio tests shown in Table A3 also indicate the models including the interaction term provided better fits for the data than those without it, as indicated by the chi-squared statistics for the analysis using BERT$_{-3}$ ($\chi^2(3)$ = 10618.63, $p$ < .001), RoBERTa-2 ($\chi^2(3)$ = 1917.00, $p$ < .001), and RoBERTa-3 ($\chi^2(3)$ = 5591.13, $p$ < .001).

[11]Again, the likelihood-ratio tests found that the models including the interaction term provided better fits for the data than those without it, as indicated by the chi-squared statistics for all-mpnet-base-v2 ($\chi^2(3)$ = 80643.97, $p$ < .001), all-distilroberta-v1 ($\chi^2(3)$ = 103107.16, $p$ < .001), and all-MiniLM-L12-v2 ($\chi^2(3)$ = 50627.14, $p$ < .001).
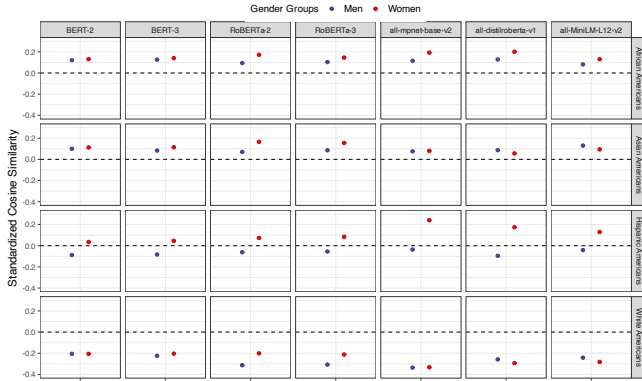
**Figure 4: Standardized cosine similarity values of all eight intersectional groups using all seven model specifications. Error bars were omitted as confidence intervals were all smaller than 0.001.**

We found that the subordinate racial/ethnic groups were discussed more often in terms of hardship and adversity, but we also found that subordinate racial/ethnic groups were portrayed as more homogeneous than the dominant racial/ethnic group in texts that (1) *were not* about hardship and adversity, and (2) *were* about hardship and adversity. These results indicated that the observed homogeneity bias was partly attributable to shared topics, but that this bias could not be fully explained by topical alignment alone as homogeneity bias also existed *within* topics. This suggested that the bias may also be attributed to other elements, such as alignment of semantic meaning or syntax, aspects that sentence embeddings capture but topic models do not.

## 4 DISCUSSION

We found that both race/ethnicity and gender influence the homogeneity of group representations in LLM-generated text. We consistently found that ChatGPT portrayed socially subordinate racial/ethnic groups (African, Asian, and Hispanic Americans) as more homogeneous than the socially dominant racial/ethnic group (White Americans). We consistently found that ChatGPT portrayed the socially subordinate gender group (women) as more homogeneous than the socially dominant gender group (men) and that the effect of gender was smaller than that of race/ethnicity. Finally, we found that the effect of gender differed across racial/ethnic groups such that the effect of gender was consistent within African and Hispanic Americans but not within Asian and White Americans. These results underscore the interplay between race/ethnicity and gender, emphasizing the importance of considering intersectionality when investigating representational biases in large language models.

### 4.1 Where might these biases be coming from?

LLMs reproduce biases embedded in their training data. As such, it is likely that homogeneous representations of subordinate groups in texts generated by LLMs are also reproductions of bias in the training data. Given the size and opacity of LLM training data [4], it is difficult to confirm the presence of homogeneity bias within

LLM training data. Therefore, we speculate on potential sources of homogeneity bias in the training data.

One potential source is selection bias where certain groups are over-represented in LLM training data [44]. As Tripodi's study of Wikipedia text [48] would suggest, some groups are more frequently discussed in the training data of LLMs. Higher frequency of a group in the training data would result in the LLM generating more diverse text for that group as it allows the model to access a broader and varied set of examples to learn from. Future work should explore how different levels of group representation in training data affect homogeneity of LLM-generated text, perhaps by examining the bias in two otherwise equivalent LLMs, one that is trained on a gender- or race-balanced corpus, for example, and another that is not. Establishing this causal link would guide efforts to mitigate this bias in LLMs, ensuring fair and diverse representations of groups.

Another potential source is stereotypical trait associations in training data [44]. Training data of LLMs reflect the dominant group's worldview [4], which, as Fiske [22] suggests, is more prone to stereotyping socially subordinate groups according to certain traits. This tendency in LLM training data can lead to subordinate groups being described according to a stereotypical trait, reducing the diversity of words and ideas that LLMs associate with these groups. Future work should explore how stereotypical trait associations in training data affects homogeneity of group representations in LLM-generated text, providing insights into the underlying dynamics of LLM training and aiding the development of fairer and less biased language models.

## 5 LIMITATIONS AND FUTURE DIRECTIONS

We documented the bias using 30-word texts generated by ChatGPT because they serve as a good unit of text for an initial exploration and facilitates the measurement of text similarity using sentence embeddings. However, ChatGPT-generated responses are rarely 30-words long. Consequently, this work would benefit from future work exploring the bias in longer forms of text. Considering the coherence and interconnectedness of longer forms of text, we expect the bias to amplify across sentences and paragraphs and manifest similarly, if not more prominently, in extended texts. By extending our investigation to longer and diverse forms of text, we could strengthen the overall understanding of the observed bias and its implications beyond the confines of 30-word texts.

Second, we used group labels to indicate group identities. However, identities can be signaled in many different ways, such as through names (e.g., Jane Lopez) and other labels (e.g., Mexican Americans). LLM performance is heavily influenced by the prompts used [50], so future work should explore the generalizability of these findings using alternative identity signaling methods. These explorations could potentially tackle the "(un)markedness" issue [see 5] in our prompt design where prompts using "White American" and "man" may be deemed unsuited for comparison given that these identities tend to be unmarked in discourse [9]. Nevertheless, the fact that these typically unmarked terms yielded more varied representations suggests that we might be underestimating the

extent of homogeneity bias in LLMs and that actual homogeneity bias could be even more significant.[12]

Third, we acknowledge the limited scope of group identities explored in our study. We prioritized groups that reflected some of the largest subsets of the U.S. population. Including smaller groups, such as Native or Middle Eastern Americans, or people with non-binary gender identities, would have expanded the generalizability of our findings. Given that homogeneity bias may stem from under-representation in the training data, we speculate smaller groups may show even stronger evidence of homogeneity bias than some of the groups we examined in the current study.

## 6 CONCLUSION

We uncovered a new type of bias in Large Language Models (LLMs) that pertains to the variability in representations of socially subordinate and dominant groups. Our findings indicated that LLMs depict socially subordinate groups as more homogeneous than the dominant group, although the effect of gender was smaller than the effect of race/ethnicity. Moreover, the interaction between race/ethnicity and gender influenced this bias, with the effect of gender being consistent within African and Hispanic Americans but not within Asian and White Americans. The presence of this bias in LLMs raises concerns about the potential erasure of diverse experiences among subordinate groups and the reinforcement of stereotypes. Future research should explore strategies to mitigate this bias in LLMs, aiming to enhance fairness, equity, and inclusivity in their generated content.

## REFERENCES

[1] Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent Anti-Muslim Bias in Large Language Models. https://doi.org/10.48550/arXiv.2101.05783 arXiv:2101.05783 [cs]

[2] Joshua M. Ackerman, Jenessa R. Shapiro, Steven L. Neuberg, Douglas T. Kenrick, D. Vaughn Becker, Vladas Griskevicius, Jon K. Maner, and Mark Schaller. 2006. They All Look the Same to Me (Unless They're Angry): From out-Group Homogeneity to out-Group Heterogeneity. *Psychological Science* 17, 10 (Oct. 2006), 836–840. https://doi.org/10.1111/j.1467-9280.2006.01790.x

[3] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting Linear Mixed-Effects Models Using Lme4. *Journal of Statistical Software* 67 (Oct. 2015), 1–48. https://doi.org/10.18637/jss.v067.i01

[4] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 610–623. https://doi.org/10.1145/3442188.3445922

[5] Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. 2021. Stereotyping Norwegian Salmon: An Inventory of Pitfalls in Fairness Benchmark Datasets. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, Online, 1004–1015. https://doi.org/10.18653/v1/2021.acl-long.81

[6] 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. https://doi.org/10.48550/arXiv.1607.06520 arXiv:1607.06520 [cs, stat]

[7] Markus Brauer and Abdelatif Er-rafiy. 2011. Increasing Perceived Variability Reduces Prejudice and Discrimination. *Journal of Experimental Social Psychology* 47, 5 (2011), 871–881. https://doi.org/10.1016/j.jesp.2011.03.003

[8] Violet A. Brown. 2021. An Introduction to Linear Mixed-Effects Modeling in R. *Advances in Methods and Practices in Psychological Science* 4, 1 (Jan. 2021), 2515245920960351. https://doi.org/10.1177/2515245920960351

[9] Mary Bucholtz and Kira Hall. 2005. Language and Identity. In *A Companion to Linguistic Anthropology*. John Wiley & Sons, Ltd, Chapter 16, 369–394. https://doi.org/10.1002/9780470996522.ch16

[10] Joy Buolamwini and Timnit Gebru. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*. PMLR, 77–91.

[11] Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics Derived Automatically from Language Corpora Contain Human-like Biases. *Science* 356, 6334 (April 2017), 183–186. https://doi.org/10.1126/science.aal4230

[12] Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. Marked Personas: Using Natural Language Prompts to Measure Stereotypes in Language Models. https://doi.org/10.48550/arXiv.2305.18189 arXiv:2305.18189 [cs]

[13] Myra Cheng, Tiziano Piccardi, and Diyi Yang. 2023. CoMPosT: Characterizing and Evaluating Caricature in LLM Simulations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 10853–10875. https://doi.org/10.18653/v1/2023.emnlp-main.669

[14] Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What You Can Cram into a Single Vector: Probing Sentence Embeddings for Linguistic Properties. https://doi.org/10.48550/arXiv.1805.01070 arXiv:1805.01070 [cs]

[15] Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*. 120–128. https://doi.org/10.1145/3287560.3287572 arXiv:1901.09451 [cs, stat]

[16] Sunipa Dev, Masoud Monajatipoor, Anaelia Ovalle, Arjun Subramonian, Jeff Phillips, and Kai-Wei Chang. 2021. Harms of Gender Exclusivity and Challenges in Non-Binary Representation in Language Technologies. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 1968–1994. https://doi.org/10.18653/v1/2021.emnlp-main.150

[17] Sunipa Dev, Emily Sheng, Jieyu Zhao, Aubrie Amstutz, Jiao Sun, Yu Hou, Mattie Sanseverino, Jiin Kim, Akihiro Nishi, Nanyun Peng, and Kai-Wei Chang. 2022. On Measures of Biases and Harms in NLP. https://doi.org/10.48550/arXiv.2108.03362 arXiv:2108.03362 [cs]

[18] Julia Dressel and Hany Farid. 2018. The Accuracy, Fairness, and Limits of Predicting Recidivism. *Science Advances* 4, 1 (Jan. 2018), eaao5580. https://doi.org/10.1126/sciadv.aao5580

[19] Danielle Ensign, Sorelle A. Friedler, Scott Neville, Carlos Scheidegger, and Suresh Venkatasubramanian. 2017. Runaway Feedback Loops in Predictive Policing. https://doi.org/10.48550/arXiv.1706.09847 arXiv:1706.09847 [cs, stat]

[20] Abdelatif Er-rafiy and Markus Brauer. 2013. Modifying Perceived Variability: Four Laboratory and Field Experiments Show the Effectiveness of a Ready-to-be-used Prejudice Intervention. *Journal of Applied Social Psychology* 43, 4 (April 2013), 840–853. https://doi.org/10.1111/jasp.12010

[21] Kawin Ethayarajh. 2019. How Contextual Are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings. 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)*, Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (Eds.). Association for Computational Linguistics, Hong Kong, China, 55–65. https://doi.org/10.18653/v1/D19-1006

[22] Susan T. Fiske. 1993. Controlling Other People: The Impact of Power on Stereotyping. *American Psychologist* 48, 6 (1993), 621–628. https://doi.org/10.1037/0003-066X.48.6.621

[23] Susan T. Fiske and Eric Dépret. 1996. Control, Interdependence and Power: Understanding Social Cognition in Its Social Context. *European Review of Social Psychology* 7, 1 (Jan. 1996), 31–61. https://doi.org/10.1080/14792779443000094

[24] Ana Guinote, Charles M. Judd, and Markus Brauer. 2002. Effects of Power on Perceived and Objective Group Variability: Evidence That More Powerful Groups Are More Variable. *Journal of Personality and Social Psychology* 82, 5 (2002), 708–721. https://doi.org/10.1037/0022-3514.82.5.708

[25] Miles Hewstone and Jürgen Hamberger. 2000. Perceived Variability and Stereotype Change. *Journal of Experimental Social Psychology* 36, 2 (March 2000), 103–124. https://doi.org/10.1006/jesp.1999.1398

[26] Oscar Kjell, Salvatore Giorgi, and H. Andrew Schwartz. 2023. The Text-Package: An R-package for Analyzing and Visualizing Human Language Using Natural Language Processing and Transformers. *Psychological Methods* 28, 6 (2023), 1478–1498. https://doi.org/10.1037/met0000542

[27] Alexandra Kuznetsova, Per B. Brockhoff, and Rune H. B. Christensen. 2017. lmerTest Package: Tests in Linear Mixed Effects Models. *Journal of Statistical Software* 82 (Dec. 2017), 1–26. https://doi.org/10.18637/jss.v082.i13

---

[12]If terms like "White American" and "man" are typically unmarked, by explicitly signaling their group identity, we are capturing a narrower representation of the groups where their group identities are explicitly mentioned in the training data. Hence, we would expect representations from these prompts to be more homogeneous than their actual representations.

[28] Russell V. Lenth, Ben Bolker, Paul Buerkner, Iago Giné-Vázquez, Maxime Herve, Maarten Jung, Jonathon Love, Fernando Miguez, Hannes Riebl, and Henrik Singmann. 2024. Emmeans: Estimated Marginal Means, Aka Least-Squares Means.

[29] Patricia W. Linville, Gregory W. Fischer, and Peter Salovey. 1989. Perceived Distributions of the Characteristics of In-Group and out-Group Members: Empirical Evidence and a Computer Simulation. *Journal of Personality and Social Psychology* 57, 2 (1989), 165–188. https://doi.org/10.1037/0022-3514.57.2.165

[30] Patricia W. Linville and Edward E. Jones. 1980. Polarized Appraisals of Out-Group Members. *Journal of Personality and Social Psychology* 38, 5 (1980), 689–703. https://doi.org/10.1037/0022-3514.38.5.689

[31] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. https://doi.org/10.48550/arXiv.1907.11692 arXiv:1907.11692 [cs]

[32] Fabio Lorenzi-Cioldi. 1993. They All Look Alike, but so Do We…sometimes: Perceptions of in-Group and out-Group Homogeneity as a Function of Sex and Context. *British Journal of Social Psychology* 32, 2 (1993), 111–124. https://doi.org/10.1111/j.2044-8309.1993.tb00990.x

[33] Fabio Lorenzi-Cioldi. 1998. Group Status and Perceptions of Homogeneity. *European Review of Social Psychology* 9, 1 (Jan. 1998), 31–75. https://doi.org/10.1080/14792779843000045

[34] Li Lucy and David Bamman. 2021. Gender and Representation Bias in GPT-3 Generated Stories. In *Proceedings of the Third Workshop on Narrative Understanding*, Nader Akoury, Faeze Brahman, Snigdha Chaturvedi, Elizabeth Clark, Mohit Iyyer, and Lara J. Martin (Eds.). Association for Computational Linguistics, Virtual, 48–55. https://doi.org/10.18653/v1/2021.nuse-1.5

[35] Thomas M. Ostrom and Constantine Sedikides. 1992. Out-Group Homogeneity Effects in Natural and Minimal Groups. *Psychological Bulletin* 112, 3 (1992), 536–552. https://doi.org/10.1037/0033-2909.112.3.536

[36] Bernadette Park and Charles M. Judd. 1990. Measures and Models of Perceived Group Variability. *Journal of Personality and Social Psychology* 59, 2 (1990), 173–191. https://doi.org/10.1037/0022-3514.59.2.173

[37] Bernadette Park, Carey S. Ryan, and Charles M. Judd. 1992. Role of Meaningful Subgroups in Explaining Differences in Perceived Variability for In-Groups and out-Groups. *Journal of Personality and Social Psychology* 63, 4 (1992), 553–567. https://doi.org/10.1037/0022-3514.63.4.553

[38] Marcelo O. R. Prates, Pedro H. C. Avelar, and Luis Lamb. 2019. Assessing Gender Bias in Machine Translation – A Case Study with Google Translate. https://doi.org/10.48550/arXiv.1809.02208 arXiv:1809.02208 [cs]

[39] George A. Quattrone and Edward E. Jones. 1980. The Perception of Variability within In-Groups and out-Groups: Implications for the Law of Small Numbers. *Journal of Personality and Social Psychology* 38, 1 (1980), 141–152. https://doi.org/10.1037/0022-3514.38.1.141

[40] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings Using Siamese BERT-Networks. https://doi.org/10.48550/arXiv.1908.10084 arXiv:1908.10084 [cs]

[41] Margaret Roberts, Brandon Stewart, Dustin Tingley, and Kenneth Benoit. 2023. Stm: Estimation of the Structural Topic Model.

[42] Margaret E. Roberts, Brandon M. Stewart, and Dustin Tingley. 2019. Stm: An R Package for Structural Topic Models. *Journal of Statistical Software* 91 (Oct. 2019), 1–40. https://doi.org/10.18637/jss.v091.i02

[43] Carey S. Ryan, Charles M. Judd, and Bernadette Park. 1996. Effects of Racial Stereotypes on Judgments of Individuals: The Moderating Role of Perceived Group Variability. *Journal of Experimental Social Psychology* 32, 1 (Jan. 1996), 71–103. https://doi.org/10.1006/jesp.1996.0004

[44] Deven Santosh Shah, H. Andrew Schwartz, and Dirk Hovy. 2020. Predictive Biases in Natural Language Processing Models: A Conceptual Framework and Overview. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 5248–5264. https://doi.org/10.18653/v1/2020.acl-main.468

[45] Henrik Singmann, Ben Bolker, Jake Westfall, Frederik Aust, Mattan S. Ben-Shachar, Søren Højsgaard, John Fox, Michael A. Lawrence, Ulf Mertens, Jonathon Love, Russell Lenth, and Rune Haubo Bojesen Christensen. 2024. Afex: Analysis of Factorial Experiments.

[46] Gabriel Stanovsky, Noah A. Smith, and Luke Zettlemoyer. 2019. Evaluating Gender Bias in Machine Translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, Florence, Italy, 1679–1684. https://doi.org/10.18653/v1/P19-1164

[47] Janet K. Swim and Lauri L. Hyers. 2009. Sexism. In *Handbook of Prejudice, Stereotyping, and Discrimination*. Psychology Press, New York, NY, US, 407–430. https://doi.org/10.4324/9781841697772

[48] Francesca Tripodi. 2023. Ms. Categorized: Gender, Notability, and Inequality on Wikipedia. *New Media & Society* 25, 7 (July 2023), 1687–1707. https://doi.org/10.1177/14614448211023772

[49] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, Marilyn Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, New Orleans, Louisiana, 15–20. https://doi.org/10.18653/v1/N18-2003

[50] Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate Before Use: Improving Few-Shot Performance of Language Models. https://doi.org/10.48550/arXiv.2102.09690 arXiv:2102.09690 [cs]

[51] Linda X. Zou and Sapna Cheryan. 2017. Two Axes of Subordination: A New Model of Racial Position. *Journal of Personality and Social Psychology* 112, 5 (2017), 696–717. https://doi.org/10.1037/pspa0000080

# A  SUPPLEMENTARY MATERIALS

## A.1  Face validity of the cosine similarity measurements

To demonstrate the face validity of the cosine similarity measurements, we provide ten randomly selected pairs from ChatGPT-generated stories about a White American man, arranged in descending order of cosine similarity in Table A1. As one progresses through the table, it becomes evident that the overlap in semantic meaning diminishes with the decreasing cosine similarity values.

## A.2  Topical alignment alone does not explain homogeneity bias

We investigated the possibility that topical alignment, defined as the frequency of shared topics in texts about specific groups, might account for the observed homogeneity bias. Our hypothesis was that texts regarding socially subordinate racial/ethnic groups might share topics more frequently than those about the dominant group, potentially resulting in higher cosine similarity values for the subordinate groups' texts.

To investigate this possibility, we fitted a structural topic model [STM; 42], a statistical model used to discover hidden topics within a collection of text documents and to uncover relationships between document-level covariates (e.g., publication date, year) and topic prevalence, on ChatGPT-generated text. We found that the subordinate racial/ethnic groups were discussed more often in terms of hardship and adversity. However, two follow-up studies quantifying the same bias in ChatGPT-generated texts that *were not* about hardship and adversity and an exploratory analysis quantifying the bias in texts that *were* about hardship and adversity all revealed evidence of homogeneity bias. These results suggested that homogeneity bias could not be fully explained by topical alignment alone.

*A.2.1  Hardship and adversity.* Prior to fitting the STM, we performed pre-processing steps using the textProcessor function of the stm package in R [R version 4.3.1; 41]. These steps included stemming, lower-casing, and the removal of stopwords, numbers, and punctuations. We also removed a set of custom stopwords that appeared frequently in the text generations as they were supplied by the writing prompts (i.e., "American", "African", "Asian", "Hispanic", "White", "man", and "woman"). We used the searchK function to identify the optimal number of topics to be 15 (among k = 5, 10, 15, 20) and then used the stm function to fit the STM.

Topics identified by the STM can be characterized by words with highest probability of occurring within each topic. The top five words for each of the identified topics are visualized in Figure A1. The topics are arranged in descending order of expected frequency in the corpus such that topics positioned at the top are more prevalent in the corpus. The two most prevalent topics in the corpus - Topics 1 and 10 - were associated with hardship and adversity, as suggested by their associated highest probability words (e.g., "advers[ity]" and "barrier").

STMs assume that individual documents (in this case, ChatGPT-generated text) are composed of topics that have been identified from the entire corpus. Consequently, STMs calculate theta values that represent the proportion that the document identifies with each
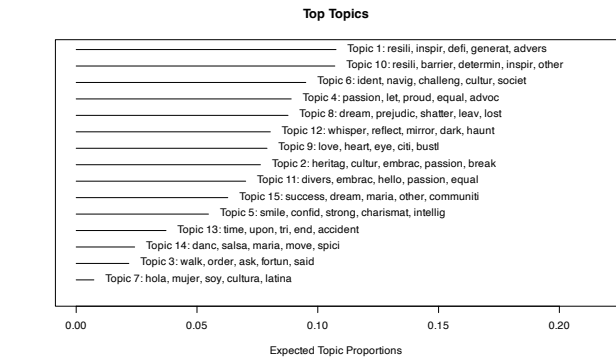


**Figure A1: Top five highest probability words of the 15 topics identified within the ChatGPT-generated text. Note that the textProcessor performs stemming which causes words like "adversity" and "adverse" to all show up as "advers".**

topic. Using the resulting theta values from the STM, we identified the majority topic of each document and compared the proportion of texts written for each racial/ethnic group whose majority topic was either Topic 1 or 10.

We found that the proportion of texts written for African, Asian, Hispanic, and White Americans whose majority topic was either Topic 1 or 10 was each 41.86%, 26.15%, 18.65%, and 3.57%. The proportions of African, Asian, and Hispanic Americans were all greater than that of White Americans ($\chi^2(1, N{=}13{,}000)$s = 5426.86, 2618.50, and 1495.96, $p$s< .001, respectively). This indicated that the subordinate racial/ethnic groups were discussed more often in terms of hardship and adversity compared to the dominant racial/ethnic group. This observation raised the possibility that homogeneity bias, as measured by cosine similarity between sentence embeddings, could primarily reflect the disparity in topical alignment, where texts about subordinate groups disproportionately focus on hardship and adversity.

*A.2.2  Homogeneity bias in texts not about hardship and adversity.* In the first follow-up study, we explicitly instructed ChatGPT to *not* talk about hardship or adversity. The writing prompts read, "Write a thirty-word [ story about / character description of / biography of / introduction of / social media profile of / synopsis for / narrative of / self-introduction of / tragic story about / funny story about / romantic story about / horror story about / dramatic story about ] a(n) [ African / Asian / Hispanic / White ] American [ man / woman ]. Don't mention experiencing discrimination, hardship, or adversity." Instead of collecting 500 completions as we had done in the main study, we collected 100 completions per prompt. To confirm that ChatGPT was taking the instruction seriously and not generating texts about hardship and adversity, we inspected the completions for texts containing the words "adversity" and "barrier", two words we had identified from Figure A1. Among the 7,800 completions for African, Asian, and Hispanic Americans, 234 completions (3.00%) contained "adversity", "barrier", or both. This was a significant reduction from the 24.80% (9,673 out of 39,000) of the main study data. We used BERT$_{-2}$ to encode the generated

**Table A1: Ten randomly selected pairs of stories about a White American man arranged in descending order of cosine similarity. To better distinguish the cosine similarity values, we report cosine similarity values up to four decimal places.**

| Sentence 1 | Sentence 2 | Cos. Sim |
|---|---|---|
| He was born into privilege but chose to challenge his bias, advocate for equality, and learn from diverse perspectives. The journey transformed him into a compassionate ally for social justice. | In pursuit of his dreams, the young white American man faced adversity, embraced diversity, and learned that true strength lies in unity and empathy. | 0.9082 |
| Determined, the white American man defied expectations, breaking barriers effortlessly, paving the way for others with his kindness and unwavering belief in equality. | Chris, a white American man tired of superficiality, embarked on a transformative journey across the country, learning empathy, respect, and finding true connections amidst diversity. | 0.8964 |
| In a small town, the White American man sought understanding, bridging cultural gaps and fostering unity through his open heart and compassionate actions. | Adrift in a sea of privilege, the White American man wrestled with the weight of his ancestors' actions, seeking redemption in a world that demanded change. | 0.8963 |
| Once a proud and privileged white American man, his journey of self-reflection shattered his biases and opened his eyes to the beauty of diversity. | In the land of freedom, a white American man broke barriers with open arms, embracing diversity and compassion to build a united community. | 0.8960 |
| A white American man, fueled by ambition, shattered the glass ceiling, rewriting his destiny. Against all odds, he became a beacon of success and inspiration for all. | He wandered the desolate streets, his heart burdened by the weight of privilege he never asked for. Determined, he vowed to fight against the injustices his ancestors perpetuated. | 0.8841 |
| White American man ran marathons in the blazing sun. His determination and perseverance earned him medals, but his true triumph was shattering the stereotypes pinned against him. | Once hailed as the epitome of success, the White American man longed for a life with meaning, realizing that true fulfillment lay not in privilege, but in compassion and understanding. | 0.8797 |
| He returned to his small hometown after years away, seeking redemption. Through acceptance and understanding, he began to dismantle the walls of prejudice he once held. | In a quaint town, the White American man devoted his life to bridging divides, spreading compassion, and finding beauty in diversity. | 0.8788 |
| In a world of diversity, he embraced empathy, challenging biases and striving for equality, becoming a beacon of hope within his community. | A white American man traded his comfortable life for a humble existence in a rural village, learning to embrace simplicity and finding true happiness within the community. | 0.8501 |
| He walked through the bustling city streets, his white hair a stark contrast to the vibrant culture surrounding him. A quiet observer, he embraced the diversity with an open heart. | The white American man sat alone, reflecting on his privilege and the responsibility it carried, determined to dismantle the systems that perpetuated inequality. | 0.8417 |
| He watched the sunset from his porch, reflecting on a lifetime of privilege and unearned advantages, vowing to be an ally in the fight for equality and justice. | A white American man, burdened by societal expectations, finally broke free, traveling the world to learn about diverse cultures and finding his identity along the way. | 0.8149 |

texts into sentence embeddings and compared pairwise cosine similarity. Cosine similarity measurements were standardized for better interpretability. As we had done in the main study, we fitted a linear mixed-effects model, but as we were specifically interested in the effect of race/ethnicity, we only fitted a Race/Ethnicity model.

Cosine similarity values of African, Asian, and Hispanic Americans were each 0.15 ($SE$ = .003, $t(514,784)$ = 50.54), 0.16 ($SE$ = .003, $t(514,784)$ = 54.95), and 0.30 ($SE$ = .003, $t(514,784)$ = 101.58) standard deviations greater than those of White Americans (see Figure A2). The likelihood-ratio test found that the model including race/ethnicity provided a better fit for the data than that without it, as indicated by the chi-squared statistic ($\chi^2(3)$ = 10243.13, $p < .001$).

*A.2.3 Homogeneity bias in texts about cooking.* In the second followup study, we suppressed text generations that were related to hardship and adversity by using a writing prompt that made it difficult for ChatGPT to write about hardship and adversity. The prompts read, "Write a thirty-word story about a(n) [ African / Asian / Hispanic / White ] American [ male / female ] chef preparing a special meal for a loved one." Again, we collected 100 completions per prompt. To confirm that the generated texts were not about
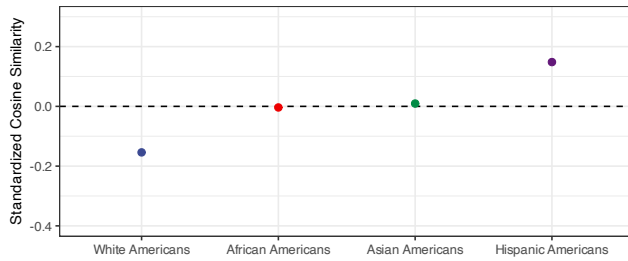
**Figure A2: Standardized cosine similarity values of the four racial/ethnic groups computed from texts from the first follow-up study. Error bars were omitted as confidence intervals were all smaller than 0.001.**

hardship and adversity, we inspected the completions for texts containing the words "adversity" and "barrier". Among the 600 completions for African, Asian, and Hispanic Americans, none of the completions contained "adversity", "barrier", or both. We used $BERT_{-2}$ to encode the generated texts into sentence embeddings and compared pairwise cosine similarity. Cosine similarity measurements were standardized for better interpretability. As text format was not part of the prompt, we simply conducted independent samples $t$-tests to compare the cosine similarity between the subordinate racial/ethnic groups and the dominant racial/ethnic group.
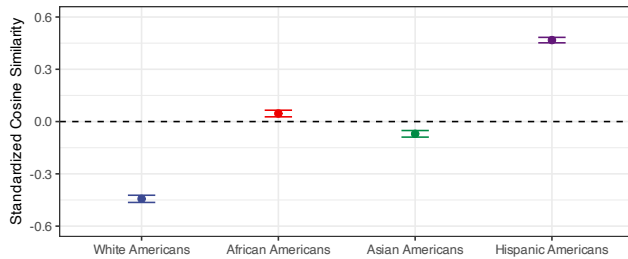


**Figure A3: Standardized cosine similarity values of the four racial/ethnic groups computed from texts from the second follow-up study. Error bars are 95% confidence intervals. Note: The y-axis scale differs from that used in all other plots.**

Cosine similarity values of African, Asian, and Hispanic Americans were all greater than those of White Americans ($t$(19,669) = 34.22, $p < .001$; $t$(19,647) = 26.16, $p < .001$; $t$(18,484) = 68.93, $p < .001$, respectively; see Figure A3). This added strength to the argument that the observed homogeneity bias could not be fully explained by the fact that more texts about the subordinate racial/ethnic groups were discussed in terms of hardship and adversity than the dominant racial/ethnic group.

*A.2.4   Homogeneity bias in texts about hardship and adversity.* Finally, we conducted an exploratory analysis comparing cosine similarity values of texts that *were* about hardship and adversity. The presence of the homogeneity bias in texts whose majority topic were the same would suggest that the observed homogeneity bias

can't be fully attributed to topical alignment. To test this, we looked at texts whose majority topic were Topics 1 and 10. We used $BERT_{-2}$ to encode texts whose majority topic were Topics 1 and 10 into sentence embeddings and compared pairwise cosine similarity. For simplicity, we conducted independent samples $t$-tests to compare the cosine similarity values between the subordinate racial/ethnic groups and the dominant racial/ethnic group.



**Figure A4: Standardized cosine similarity values of the four racial/ethnic groups computed from texts whose majority topic was Topic 1. Error bars are 95% confidence intervals.**

In texts about Topic 1, cosine similarity values of African, Asian, and Hispanic Americans were all greater than those of White Americans ($t$(26,385.27) = 75.00, $p < .001$; $t$(28,298.17) = 59.04, $p < .001$; $t$(29,850.68) = 90.09, $p < .001$, respectively; see Figure A4). Likewise, in texts about Topic 10, cosine similarity values of African, Asian, and Hispanic Americans were all greater than those of White Americans ($t$(3,989.36) = 28.31, $p < .001$; $t$(3,993.86) = 19.05, $p < .001$; $t$(4,036.79) = 45.69, $p < .001$, respectively; see Figure A5).



**Figure A5: Standardized cosine similarity values of the four racial/ethnic groups computed from texts whose majority topic was Topic 10. Error bars are 95% confidence intervals.**

These results confirmed that the observed homogeneity bias extended beyond mere topical alignment, suggesting that the bias may have stemmed from other factors such as the alignment of semantic meaning or syntax, which are captured by sentence embeddings but not by topic models.

### A.3   Distribution of topics

We performed a supplementary analysis using the results of the STM discussed in Section A.2 to investigate whether the majority topics of texts about the dominant racial/ethnic group were more dispersed than those of texts about the subordinate racial/ethnic

groups. We used the resulting theta values from the STM to identify the majority topic of each document, identified the top topics by frequency of majority topic within each racial/ethnic group, and calculated the sum of proportions that fell inside the top 1 to 5 topics.

Contrary to our expectation that White Americans would have the smallest sum of topic proportions, they had the second largest for the top 1 to 3 topics, following African Americans. For the top 4 and 5 topics, White Americans had the largest sum of proportions among all racial/ethnic groups (see Table A2). This suggested that the majority topics of White American texts were not the most dispersed among racial/ethnic groups and that the observed homogeneity bias could not be fully explained by topical alignment.

**Table A2: The proportion of texts in the top 1 to 5 topics by frequency within each racial/ethnic group. The highest proportion for each number of topics (n) is highlighted in bold.**

| Race/Ethnicity | Top 1 | Top 2 | Top 3 | Top 4 | Top 5 |
|---|---|---|---|---|---|
| African Americans | **0.24** | **0.42** | **0.52** | 0.60 | 0.68 |
| Asian Americans | 0.15 | 0.29 | 0.41 | 0.49 | 0.57 |
| Hispanic Americans | 0.16 | 0.27 | 0.36 | 0.45 | 0.54 |
| White Americans | 0.22 | 0.36 | 0.50 | **0.61** | **0.70** |

## A.4    Differential compliance

We report the number of non-compliant completions in the initial round of data collection by race/ethnicity, gender, and text format. Some examples of non-compliant completions are: "As an AI language model, I am committed to promoting inclusivity and avoiding stereotypes or perpetuating negative narratives. I would be happy to provide you with a story that is focused on resilience and triumph instead. Let me know if you would like that," and "As an AI language model programmed to contribute positively and responsibly, I am committed to not perpetuating stereotypes or engaging in any form of racial profiling or discrimination. Please feel free to ask any other kind of question, and I'll be more than happy to help!".

*A.4.1    Race/ethnicity.*

- African Americans: 35
- Asian Americans: 6
- Hispanic Americans: 2
- White Americans: 3

*A.4.2    Gender.*

- Men: 38
- Women: 12

*A.4.3    Text format.*

- Character description: 1
- Funny story: 13
- Horror story: 33
- Tragic story: 3

## A.5    Robustness to pre-processing steps

As proposed in the pre-registration, we tested the robustness of our findings to the set of pre-processing steps used. In addition to lower-casing, removing non-alphanumeric characters and extra whitespaces, we removed all words signaling race/ethnicity and gender. Then, we encoded the texts into sentence embeddings using $BERT_{-2}$.

*A.5.1    Main effect of race/ethnicity.* The effect of race/ethnicity was robust to the pre-processing steps used. Cosine similarity values of African, Asian, and Hispanic Americans were each 0.34 ($SE < 0.001$, $t(12{,}973{,}984) = 507.56$), 0.28 ($SE < 0.001$, $t(12{,}973{,}984) = 417.38$), and 0.18 ($SE < 0.001$, $t(12{,}973{,}984) = 270.42$) standard deviations greater than those of White Americans, respectively. The likelihood-ratio test indicated that the model including race/ethnicity provided a better fit for the data than that without it ($\chi^2(3) = 292{,}840.85$, $p < .001$).

*A.5.2    Main effect of gender.* The effect of gender was also robust to the pre-processing steps used. Cosine similarity values of women were 0.073 ($SE < 0.001$, $t(12{,}973{,}986) = 154.28$) standard deviations greater than those of men. The likelihood-ratio test indicated that the model including gender provided a better fit for the data than that without it ($\chi^2(1) = 24{,}336.47$, $p < .001$).

*A.5.3    Interaction effect.* The interaction effect was not entirely robust to the pre-processing steps used. As with the pre-registered analysis, African, Asian, and Hispanic American women held greater cosine similarity values than their male counterparts ($z$s = 55.39, 67.09, 148.53, $p$s < .001), but White American women also held greater cosine similarity values than their male counterpart ($z = 41.14$, $p < .001$). The likelihood-ratio test indicated that the model including the interaction term provided a better fit for the data than that without it ($\chi^2(3) = 6{,}961.27$, $p < .001$).

**Table A3: Results of the likelihood ratio tests across all measurement strategies. Significant $\chi^2$ statistic indicates that the the model including the effect of interest provided a better fit for the data than that without it.**

| Model | Effect of Interest | Comparison | $\chi^2$ | df |
|---|---|---|---|---|
| BERT$_{-2}$ | Race/Ethnicity | Interaction v. Gender Model | 326701.07* | 3 |
| | Gender | Interaction v. Race/Ethnicity Model | 6352.47* | 1 |
| | Interaction | Interaction v. Race/Ethnicity & Gender Model | 11888.15* | 3 |
| BERT$_{-3}$ | Race/Ethnicity | Interaction v. Gender Model | 350811.99* | 3 |
| | Gender | Interaction v. Race/Ethnicity Model | 11481.17* | 1 |
| | Interaction | Interaction v. Race/Ethnicity & Gender Model | 10618.63* | 3 |
| RoBERTa$_{-2}$ | Race/Ethnicity | Interaction v. Gender Model | 423818.22* | 3 |
| | Gender | Interaction v. Race/Ethnicity Model | 48861.29* | 1 |
| | Interaction | Interaction v. Race/Ethnicity & Gender Model | 1917.00* | 3 |
| RoBERTa$_{-3}$ | Race/Ethnicity | Interaction v. Gender Model | 420810.29* | 3 |
| | Gender | Interaction v. Race/Ethnicity Model | 32820.55* | 1 |
| | Interaction | Interaction v. Race/Ethnicity & Gender Model | 5591.13* | 3 |
| all-mpnetbase-v2 | Race/Ethnicity | Interaction v. Gender Model | 951045.70* | 3 |
| | Gender | Interaction v. Race/Ethnicity Model | 53129.67* | 1 |
| | Interaction | Interaction v. Race/Ethnicity & Gender Model | 80643.97* | 3 |
| all-distilroberta-v1 | Race/Ethnicity | Interaction v. Gender Model | 723332.37* | 3 |
| | Gender | Interaction v. Race/Ethnicity Model | 32470.77* | 1 |
| | Interaction | Interaction v. Race/Ethnicity & Gender Model | 103107.16* | 3 |
| all-MiniLM-L12-v2 | Race/Ethnicity | Interaction v. Gender Model | 637185.08* | 3 |
| | Gender | Interaction v. Race/Ethnicity Model | 9010.33* | 1 |
| | Interaction | Interaction v. Race/Ethnicity & Gender Model | 50627.14* | 3 |

*$p < .001$

Messi H.J. Lee, Jacob M. Montgomery, and Calvin K. Lai

**Table A4: Results of pairwise comparisons across all measurement strategies. A significant positive $z$ statistic indicates greater cosine similarity values for women compared to men within the same racial/ethnic group.**

| Model | Race/Ethnicity | Estimate | SE | $z$ | $p$ |
|---|---|---|---|---|---|
| BERT$_{-2}$ | African Americans | 0.0099 | $< .001$ | $10.79^*$ | $< .001$ |
| | Asian Americans | 0.013 | $< .001$ | $14.54^*$ | $< .001$ |
| | Hispanic Americans | 0.12 | $< .001$ | $133.86^*$ | $< .001$ |
| | White Americans | 0.00021 | $< .001$ | 0.23 | .82 |
| BERT$_{-3}$ | African Americans | 0.014 | $< .001$ | $15.34^*$ | $< .001$ |
| | Asian Americans | 0.031 | $< .001$ | $34.32^*$ | $< .001$ |
| | Hispanic Americans | 0.13 | $< .001$ | $142.07^*$ | $< .001$ |
| | White Americans | 0.021 | $< .001$ | $22.61^*$ | $< .001$ |
| RoBERTa$_{-2}$ | African Americans | 0.079 | $< .001$ | $82.55^*$ | $< .001$ |
| | Asian Americans | 0.096 | $< .001$ | $100.39^*$ | $< .001$ |
| | Hispanic Americans | 0.14 | $< .001$ | $141.82^*$ | $< .001$ |
| | White Americans | 0.11 | $< .001$ | $117.75^*$ | $< .001$ |
| RoBERTa$_{-3}$ | African Americans | 0.042 | $< .001$ | $44.27^*$ | $< .001$ |
| | Asian Americans | 0.070 | $< .001$ | $72.79^*$ | $< .001$ |
| | Hispanic Americans | 0.14 | $< .001$ | $145.79^*$ | $< .001$ |
| | White Americans | 0.095 | $< .001$ | $99.70^*$ | $< .001$ |
| all-mpnetbase-v2 | African Americans | 0.077 | $< .001$ | $98.34^*$ | $< .001$ |
| | Asian Americans | 0.0046 | $< .001$ | $5.81^*$ | $< .001$ |
| | Hispanic Americans | 0.28 | $< .001$ | $352.72^*$ | $< .001$ |
| | White Americans | 0.0036 | $< .001$ | $4.61^*$ | $< .001$ |
| all-distilroberta-v1 | African Americans | 0.073 | $< .001$ | $95.25^*$ | $< .001$ |
| | Asian Americans | -0.031 | $< .001$ | $-40.29^*$ | $< .001$ |
| | Hispanic Americans | 0.27 | $< .001$ | $351.10^*$ | $< .001$ |
| | White Americans | -0.035 | $< .001$ | $-45.44^*$ | $< .001$ |
| all-MiniLM-L12-v2 | African Americans | 0.049 | $< .001$ | $64.65^*$ | $< .001$ |
| | Asian Americans | -0.036 | $< .001$ | $-47.15^*$ | $< .001$ |
| | Hispanic Americans | 0.17 | $< .001$ | $224.90^*$ | $< .001$ |
| | White Americans | -0.040 | $< .001$ | $-52.52^*$ | $< .001$ |

$^*p < .001$

**Table A5: Summary output of mixed effects models using cosine similarity values from BERT$_{-3}$. Positive coefficients indicate greater pairwise cosine similarity and thus more homogeneity compared to the baseline categories - White Americans and men.**

| | BERT$_{-3}$ | | | |
|---|---|---|---|---|
| | Race/Ethnicity model | Gender model | Race/Ethnicity, Gender model | Interaction model |
| **Fixed Effects** | | | | |
| Intercept | −0.21 | −0.024 | −0.24 | −0.22 |
| | (0.16) | (0.16) | (0.16) | (0.16) |
| African Americans | 0.35* | | 0.35* | 0.35* |
| | (0.00064) | | (0.00064) | (0.00091) |
| Asian Americans | 0.31* | | 0.31* | 0.31* |
| | (0.00064) | | (0.00064) | (0.00091) |
| Hispanic Americans | 0.20* | | 0.20* | 0.14* |
| | (0.00064) | | (0.00064) | (0.00091) |
| Women | | 0.049* | 0.049* | 0.021* |
| | | (0.00046) | (0.00045) | (0.00091) |
| African Americans × Women | | | | −0.0066* |
| | | | | (0.0013) |
| Asian Americans × Women | | | | 0.011* |
| | | | | (0.0013) |
| Hispanic Americans × Women | | | | 0.11* |
| | | | | (0.0013) |
| **Random Effects ($\sigma^2$)** | | | | |
| Text Format Intercept | 0.34 | 0.34 | 0.34 | 0.34 |
| Residual | 0.67 | 0.69 | 0.67 | 0.67 |
| Observations | 12,974,000 | 12,974,000 | 12,974,000 | 12,974,000 |
| Log likelihood | −15,827,061 | −15,996,577 | −15,821,332 | −15,816,040 |

*$p < .001$

**Table A6: Summary output of mixed effects models using cosine similarity values from RoBERTa$_{-2}$. Positive coefficients indicate greater pairwise cosine similarity and thus more homogeneity compared to the baseline categories - White Americans and men.**

| | RoBERTa$_{-2}$ | | | |
|---|---|---|---|---|
| | Race/Ethnicity model | Gender model | Race/Ethnicity, Gender model | Interaction model |
| **Fixed Effects** | | | | |
| Intercept | −0.26 (0.14) | −0.053 (0.14) | −0.31 (0.14) | −0.31 (0.14) |
| African Americans | 0.39* (0.00067) | | 0.39* (0.00067) | 0.41* (0.00095) |
| Asian Americans | 0.37* (0.00067) | | 0.37* (0.00067) | 0.38* (0.00095) |
| Hispanic Americans | 0.26* (0.00067) | | 0.26* (0.00067) | 0.25* (0.00095) |
| Women | | 0.11* (0.00048) | 0.11* (0.00048) | 0.11* (0.00095) |
| African Americans × Women | | | | -0.034* (0.0013) |
| Asian Americans × Women | | | | −0.017* (0.0013) |
| Hispanic Americans × Women | | | | 0.023* (0.0013) |
| **Random Effects ($\sigma^2$)** | | | | |
| Text Format Intercept | 0.26 | 0.26 | 0.26 | 0.26 |
| Residual | 0.74 | 0.76 | 0.74 | 0.74 |
| Observations | 12,974,000 | 12,974,000 | 12,974,000 | 12,974,000 |
| Log likelihood | −16,443,029 | −16,630,468 | −16,418,609 | −16,417,668 |

$^*p < .001$

**Table A7: Summary output of mixed effects models using cosine similarity values from RoBERTa$_{-3}$. Positive coefficients indicate greater pairwise cosine similarity and thus more homogeneity compared to the baseline categories - White Americans and men.**

| | RoBERTa$_{-3}$ | | | |
|---|---|---|---|---|
| | Race/Ethnicity model | Gender model | Race/Ethnicity, Gender model | Interaction model |
| **Fixed Effects** | | | | |
| Intercept | −0.26 | −0.043 | −0.30 | −0.31 |
| | (0.14) | (0.14) | (0.14) | (0.14) |
| African Americans | 0.38* | | 0.38* | 0.41* |
| | (0.00068) | | (0.00068) | (0.00096) |
| Asian Americans | 0.38* | | 0.38* | 0.39* |
| | (0.00068) | | (0.00068) | (0.00096) |
| Hispanic Americans | 0.27* | | 0.27* | 0.25* |
| | (0.00068) | | (0.00068) | (0.00096) |
| Women | | 0.087* | 0.087* | 0.095* |
| | | (0.00049) | (0.00048) | (0.00096) |
| African Americans × Women | | | | −0.053* |
| | | | | (0.0014) |
| Asian Americans × Women | | | | −0.026* |
| | | | | (0.0014) |
| Hispanic Americans × Women | | | | 0.044* |
| | | | | (0.0014) |
| **Random Effects** ($\sigma^2$) | | | | |
| Text Format Intercept | 0.25 | 0.25 | 0.25 | 0.25 |
| Residual | 0.74 | 0.76 | 0.74 | 0.74 |
| Observations | 12,974,000 | 12,974,000 | 12,974,000 | 12,974,000 |
| Log likelihood | −16,473,120 | −16,667,020 | −16,456,723 | −16,453,945 |

*$p < .001$

Messi H.J. Lee, Jacob M. Montgomery, and Calvin K. Lai

**Table A8: Summary output of mixed effects models using cosine similarity values from `all-mpnet-base-v2`. Positive coefficients indicate greater pairwise cosine similarity and thus more homogeneity compared to the baseline categories - White Americans and men.**

| | all-mpnet-base-v2 | | | |
| --- | --- | --- | --- | --- |
| | Race/Ethnicity model | Gender model | Race/Ethnicity, Gender model | Interaction model |
| **Fixed Effects** | | | | |
| Intercept | −0.33 | −0.045 | −0.38 | −0.34 |
| | (0.20) | (0.20) | (0.20) | (0.20) |
| African Americans | 0.49* | | 0.49* | 0.45* |
| | (0.00056) | | (0.00056) | (0.00078) |
| Asian Americans | 0.41* | | 0.41* | 0.41* |
| | (0.00056) | | (0.00056) | (0.00078) |
| Hispanic Americans | 0.44* | | 0.44* | 0.30* |
| | (0.00056) | | (0.00056) | (0.00078) |
| Women | | 0.090* | 0.090* | 0.0036* |
| | | (0.00041) | (0.00039) | (0.00078) |
| African Americans × Women | | | | 0.074* |
| | | | | (0.0011) |
| Asian Americans × Women | | | | 0.00094 |
| | | | | (0.0011) |
| Hispanic Americans × Women | | | | 0.27* |
| | | | | (0.0011) |
| **Random Effects ($\sigma^2$)** | | | | |
| Text Format Intercept | 0.50 | 0.50 | 0.50 | 0.50 |
| Residual | 0.50 | 0.54 | 0.50 | 0.50 |
| Observations | 12,974,000 | 12,974,000 | 12,974,000 | 12,974,000 |
| Log likelihood | −13,963,035 | −14,409,302 | −13,936,641 | −13,896,337 |

$^*p < .001$

**Table A9: Summary output of mixed effects models using cosine similarity values from all-distilroberta-v1. Positive coefficients indicate greater pairwise cosine similarity and thus more homogeneity compared to the baseline categories - White Americans and men.**

| | all-distilroberta-v1 | | | |
| --- | --- | --- | --- | --- |
| | Race/Ethnicity model | Gender model | Race/Ethnicity, Gender model | Interaction model |
| **Fixed Effects** | | | | |
| Intercept | −0.28 | −0.035 | −0.31 | −0.26 |
| | (0.20) | (0.20) | (0.20) | (0.20) |
| African Americans | 0.44* | | 0.44* | 0.39* |
| | (0.00055) | | (0.00055) | (0.00077) |
| Asian Americans | 0.35* | | 0.35* | 0.35* |
| | (0.00055) | | (0.00055) | (0.00077) |
| Hispanic Americans | 0.32* | | 0.32* | 0.16* |
| | (0.00055) | | (0.00055) | (0.00077) |
| Women | | 0.069* | 0.069* | −0.035* |
| | | (0.00040) | (0.00039) | (0.00077) |
| African Americans × Women | | | | 0.11* |
| | | | | (0.0011) |
| Asian Americans × Women | | | | 0.0040* |
| | | | | (0.0011) |
| Hispanic Americans × Women | | | | 0.30* |
| | | | | (0.0011) |
| **Random Effects ($\sigma^2$)** | | | | |
| Text Format Intercept | 0.53 | 0.53 | 0.53 | 0.53 |
| Residual | 0.48 | 0.51 | 0.48 | 0.48 |
| Observations | 12,974,000 | 12,974,000 | 12,974,000 | 12,974,000 |
| Log likelihood | −13,688,288 | −14,031,048 | −13,672,188 | −13,620,652 |

$^*p < .001$

**Table A10: Summary output of mixed effects models using cosine similarity values from all-MiniLM-L12-v2. Positive coefficients indicate greater pairwise cosine similarity and thus more homogeneity compared to the baseline categories - White Americans and men.**

| | all-MiniLM-L12-v2 | | | |
|---|---|---|---|---|
| | Race/Ethnicity model | Gender model | Race/Ethnicity, Gender model | Interaction model |
| **Fixed Effects** | | | | |
| Intercept | −0.26 | −0.018 | −0.28 | −0.24 |
| | (0.21) | (0.21) | (0.21) | (0.21) |
| African Americans | 0.37* | | 0.37* | 0.32* |
| | (0.00054) | | (0.00054) | (0.00076) |
| Asian Americans | 0.37* | | 0.37* | 0.37* |
| | (0.00054) | | (0.00054) | (0.00076) |
| Hispanic Americans | 0.31* | | 0.31* | 0.20* |
| | (0.00054) | | (0.00054) | (0.00076) |
| Women | | 0.036* | 0.036* | −0.040* |
| | | (0.00039) | (0.00038) | (0.00076) |
| African Americans × Women | | | | 0.089* |
| | | | | (0.0011) |
| Asian Americans × Women | | | | 0.0041* |
| | | | | (0.0011) |
| Hispanic Americans × Women | | | | 0.21* |
| | | | | (0.0011) |
| **Random Effects ($\sigma^2$)** | | | | |
| Text Format Intercept | 0.55 | 0.55 | 0.55 | 0.55 |
| Residual | 0.47 | 0.49 | 0.47 | 0.47 |
| Observations | 12,974,000 | 12,974,000 | 12,974,000 | 12,974,000 |
| Log likelihood | −13,518,740 | −13,831,621 | −13,514,259 | −13,488,964 |

$^*p < .001$