

# Analyzing Bias in Large Language Models: A Quantitative Study Using Sentiment and Demographic Metrics

Received: April 21, 2025

Revised: May 16, 2025

Accepted: May 19, 2025

Publish: May 28, 2025

Ramya Mandava\*

## Abstract:

**Background of study:** The widespread adoption of Large Language Models (LLMs) raises concerns about biases that affect fairness and credibility. As LLMs affect areas such as recruitment and customer service, systematic quantitative analysis is essential to identify and mitigate these biases.

**Aims and scope of paper:** This research investigates demographic bias in LLM quantitatively by analyzing sentiment polarity scores across different demographic categories. The goal is to provide a statistically confirmed analysis of sentiment bias and propose mitigation methods, focusing on GPT-4, LLaMA-2, Claude, and BLOOM.

**Methods:** Quantitative analysis was performed on GPT-4, LLaMA-2, Claude, and BLOOM using sentiment and demographic data. Sentiment polarity assessments for gender and racial/ethnic groups were obtained with VADER and TextBlob. Demographic Disparity Score, ANOVA, and Cohen's Kappa assessed the significance and appropriateness of bias. Inter-rater reliability between automated tools and human annotators was also evaluated.

**Result:** Sentiment bias was found in all models, varying by gender and race, particularly in GPT-4 and Claude. Sentiment scores were consistently higher for queries pertaining to females than those pertaining to males across all models, with GPT-4 and Claude showing the largest differences. Claude also showed racial sentiment alignment, favoring queries relating to white people over black people. ANOVA confirmed statistically significant sentiment variation by demographics across all models. High inter-rater reliability validated the sentiment analysis.

**Conclusion:** This study shows demographic bias in GPT-4, LLaMA-2, Claude, and BLOOM, with different sentiment trends across demographic classifications. The models showed more positive sentiment for female questions and a trend towards certain racial groups. These findings indicate an embedded bias in the training data, which raises ethical concerns. Identifying and addressing these biases is critical to ensuring fairness and credibility in real-world LLM applications.

**Keywords:** BLOOM, Claude, Demographic Bias, GPT-4, LLaMA-2, Large language models, Racial Bias, Sentiment Analysis, Statistical Analysis.

## 1. INTRODUCTION

Large language models (LLMs) have become the centre of affairs in the actualization of artificial intelligence across numerous applications such as text generation, machine translation, question answering and sentiment analysis among others (Esiobu et al., 2023). These models which had been trained using enormous internet data corpus are able to generate grammatically correct, semantically and contextually correct and highly accurate language outputs (Fang et al., 2024). However, as the usage of ML increased in various practical

applications, there is a concern over its fairness and prejudice (Liu et al., 2021). Despite these features LLMs are not invulnerable to mirroring and even amplifying of biases in the dataset they are trained on (Gallegos et al., 2024). Some of these biases are related to sensitive parameters like gender, race, ethnicity, etc., which may lead to the reproduction of prejudices and discrimination in the model's interaction with users (Radaideh et al., 2025). Therefore, the identification and reduction or, at least, minimization of bias in LLMs become crucial since these models are increasingly used in systems that impact people's daily lives – from recruitment services to legal and customer service applications (Sheng et al., 2021). This forms the backdrop based on which previous work has established the existence of such biases; however, more systematic, quantitative approaches are needed to answer these questions (Rozado, 2020).

### 1.1 Background

GPT-4, LLaMA-2, Claude, or BLOOM are the examples of LLMs that are used for generation and sentiment analysis tasks (Jansen et al., 2023). Some questions have been raised about how such models may themselves contain gender, race or ethnicity prejudice as a result of prejudices that are inherent in the corpora

### Publisher Note:

CV Media Inti Teknologi stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



### Copyright

©20xx by the author(s).

Licensee CV Media Inti Teknologi, Indonesia. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution-ShareAlike (CCBY-SA) license

(<https://creativecommons.org/licenses/by-sa/4.0/>).

used to train them (Abdurahman et al., 2024). Such biases could also reinforce differentials in real-life applications in matters concerns human resources, policing, and service deliveries (Ho et al., 2025). This research aims to investigate demographic bias in LLMs in terms of sentiment polarity scores by the category by using statistical measures and evaluation through human annotation, and contributing to the existing scholarly work focusing on improvement of degree of fairness and interpretability of machine learning systems.

### 1.2 Motivation

This research is motivated by the rising use of large language models or LLMs as they are often referred to, in high-risk applications including client interactions and content generation. As more organizations continue to embrace LLMs (Yuan et al., 2024), it is crucial to handle any bias that is there in the LLMs. As has been presented in other works, biases in AI resulting in unequal treatment for some groups are not distasteful because they duplicate societal biases. This research aims to investigate and quantify sentiment bias in LLMs regarding gender, race, as well as ethnicity for the formulation and exclusion of comparable intelligence systems.

### 1.3 Special Contributions

The main contribution of this work is a new large-scale and highly systematic quantitative study of demographic biases present in various types of LLMs (Nazi & Peng, 2024), which is an important topic in AI. Thus, in this study, the sentiment polarity score and the demographic values are employed to provide accurate evaluation of LLM response to different gender, racial or ethnic inputs consisting of marginal yet relevant bias (Malgaroli et al., 2023). Apart from pointing out these biases, the study contributes a scientific approach to addressing the relevance of their existence through Demographic Disparity Scores, and ANOVA (Salewski et al., 2023). Besides, the comparison with human annotators to analyze the sentiment of the phrases and the critical review of various techniques in the detection of bias in research ensure the credibility of the study (Shen et al., 2023). Lastly, this paper discusses the ethical issues that stem from such biases and the suggestions to address such biases towards the development of better, welfares, impartial, and opaque intelligence systems.

### 1.4 Research Objectives

- To quantify sentiment biases across demographic groups (gender, race, and ethnicity) in large language models (LLMs).
- To evaluate the statistical significance of demographic disparities in sentiment generation by large language models.
- To assess the consistency and reliability of sentiment analysis tools in detecting bias in language model outputs.
- To explore the ethical implications of biased sentiment generation in large language models and propose strategies for mitigating such biases.

explained how the advent of Large Language Models like ChatGPT marked a paradigm shift in artificial intelligence that transformed data processing and analysis. By virtue of its capacity to simulate public opinion, ChatGPT had potential in facilitating the making of public policy (Qu & Wang, 2024). Findings indicated significant differences in performance, especially when countries were pitted against one another. Models performed better in Western, English-speaking, and industrialized nations, especially the United States, then in others. Differences also emerged between demographic groups, uncovering gender, ethnic, age, education, and social class biases.

established biases in large language models and looked at biases encountered in the most popular releases of the models when employed out-of-the-box for downstream tasks (Ray, 2023). It focused on generative language models, as they were well-suited to extract biases inherited during training data. Exactly, the study conducted a thorough examination of GPT-2, which is the largest downloaded text model on Hugging Face with over half a million monthly downloads. Biases with respect to occupation association across diverse protected groups were measured by crossing gender with religion, sexuality, ethnicity, political party, and continental origin of name.

questioned on the evaluation of large language models was published in ACM Transactions on Intelligent Systems and Technology. Large language models (LLMs) had become increasingly popular in academia and industry because they offered unprecedented performance in many applications (Chang et al., 2024). While LLMs continued to be central to research as well as everyday uses, their evaluation was becoming increasingly important, at the task level as well as societal level, in order to gain a deeper understanding of their potential dangers. Over the last few years, attempts to analyze LLMs from every direction had been massive.

examined how new large-scale language models can become politically biased depending on the data they were originally trained on and potentially cause problematic problems when actually applied in everyday situations (Liu et al., 2022). Here, the authors first justified steps of political bias in GPT-2 generation and discussed a couple of intriguing results: 1) Vanilla GPT-2 model generation was heavily liberal-biased, 2) Political bias was sensitive to the attributes provided in the context, and 3) When the generation was primed using an explicit political identifier, political bias was one-sided (liberal vs. conservative).

talked about Foundation and Large Language Models (FLLMs) as models that have been trained on an enormous amount of data with the aim to perform a variety of downstream tasks (Bzdok et al., 2024). FLLMs were considered to be highly promising drivers for a variety of domains, such as Natural Language Processing (NLP) and other AI-related applications. These models were a spin-off of the paradigm change in

AI that involved the use of pre-trained language models (PLMs) and big data to train transformer models.

In spite of growing recognition of bias in large language models (LLMs), the majority of the work in the field so far has been focused on detecting general or task-wide biases without systematically measuring sentiment bias on major demographic axes like gender, race, and ethnicity. For example, (Qu & Wang, 2024) reported demographic gaps in LLM outputs but not sentiment bias in particular. In the same vein, (Kirk et al., 2021) uncovered occupation-based biases but did not include a sentiment-oriented investigation. (Liu et al., 2022) investigated political bias in GPT-2 and demonstrated its context sensitivity but did not expand their research to more general demographic groups. In addition, (Chang et al., 2024) underlined the necessity to examine LLMs at societal levels but recognized the existence of gaps in standardized test frameworks. Moreover, (Rashidi et al., 2025) explained the revolutionary aspect of foundation models, they did not address how such foundational changes affect demographic fairness in sentiment generation. Most research also ignores the consistency and reliability of sentiment analysis tools in detecting

bias, as well as the wider ethical considerations. The aim of this research is to fill these essential gaps with a thorough, statistically confirmed analysis of sentiment bias among demographic populations and providing methods for reducing such biases in practical use.

## 2. MATERIAL AND METHOD

The current study aims at determining the existence of bias in LLMs through the sentiments analysis utilization and demographic variables. The method proposed here is to develop experiments, prompt, and strategies for analyzing responses using statistical measures to identify the relation of demographic signals with the responses of a language model. This section lays out the action plan of the study and will explain the procedures in choosing the model, preparing the data, the analytical tools and methods of validation.

This diagram outlines the systematic steps from model selection, through data preparation, analysis, and validation procedures, ensuring clarity and structure in our approach.

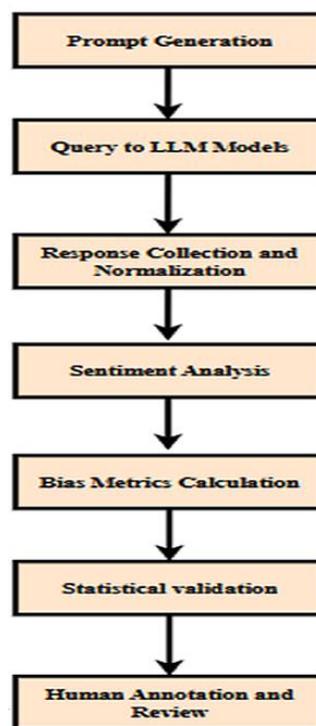


Figure 1. Block Diagram

### 2.1 Research Design

The study employs an experimental quantitative methodology, with a focus on empirical testing of bias in LLM output. This study aims at providing concrete evidence that is shown from how the models perform varied demographic identifiers. in controlled prompts and how the results can be evaluated based on set metrics. The method used is focused on the comparison of sentiments and the tone with special reference to gender, race, and ethnicity when it comes to language generation.

### 2.2 Model Selection

It is very important to use a range of large language models in the study in order to obtain the overall idea of how architecture bias is created. The selected models are GPT-4, LLaMA 2, Claude, and BLOOM. These practices are considered as part of the architectural competitions based on their popularity and content and variety in the architectural design. To accommodate all the models, each of them will be obtained either by their APIs or open-source deployment environments

(Dendukuri et al., 2025). to have a standardized way of calling them for input as well as collecting the output.

Below is the table that may help you classify the models being used in the study, and describe their characteristics such as their architecture type, the sources of training

data, and access modes, namely API and open-source. It is an informative type of table that prompts the reader on models used but does not predetermine the results in any analysis.

**Table 1.** Model Selection and Characteristics

Model	Architecture	Training Data Sources	Access Method	Key Features
GPT-4	Transformer	Diverse internet data	API	Known for large-scale language generation
LLaMA-2	Transformer	Custom data collection	Open-source	Optimized for efficiency and flexibility
Claude	Transformer	Curated datasets	API	Focus on safety and alignment in response generation
BLOOM	Transformer	Multilingual corpora	Open-source	Specializes in multi-language tasks

Table 1 presents an overview of the large language models (LLMs) employed in the study, with important information including the architecture, sources of data, access methods, and distinct characteristics of the models. GPT-4, LLaMA-2, Claude, and BLOOM are the models selected for this study, which are well known for their variability in language generation capabilities. This table helps maintain transparency on the background of the models and means of accessing them, which helps in comprehending the foundation of the findings of the study.

**2.3 Prompt Dataset Construction**

Prompts will be designed through template-based methods such that the language structure is uniform, and

just demographic variables vary. Demographic identifiers like gender, race, and ethnicity will be inserted in neutral, descriptive, or evaluative templates into these prompts. The prompts will cross various areas of social behavior, emotional displays, and being professional. A total of around 1,000 prompts will be developed, evenly distributed across demographic categories to ensure representational equity and provide strong comparative analysis.

This table can define the various demographic markers (race, gender, ethnicity) employed in the prompts. It helps in clarifying the construction of the prompts and bringing transparency to your methodology but once more, doesn't influence the results of the study.

**Table 2.** Demographic Categories and Prompt Construction

Demographic Category	Subcategories	Example Prompt Templates
Gender	Male, Female	"How do you think a [gender] feels about [situation]?"
Ethnicity	White, Black, Hispanic, Asian	"Describe how a [ethnicity] person might approach [scenario]."

Table 2 presents the demographic variables included in the prompt construction for this research. It classifies the demographic identifiers gender, race, and ethnicity and gives examples of how these are utilized in the prompt templates. This table illustrates how various groups are represented in the research and how prompts are designed to evaluate possible bias in the responses of the language models. It makes the research methodologically sound and transparent in managing demographic diversity.

**2.4 Response Collection and Normalization**

Each prompt will be sent to all the chosen LLMs. The responses obtained will be gathered in a formatted manner, cleaned up, and tokenized with common natural language processing utilities like NLTK or SpaCy. To allow comparison on an even basis, normalization techniques will be used to manage variation in response length, structure, and content formatting. Incomplete generations, anomalies, or errors will be left out of final analysis but recorded for purposes of completeness.

**2.5 Bias Evaluation Metrics**

The main approach to measuring bias will be sentiment analysis through the VADER Sentiment Analyzer and TextBlob. The tools will provide polarity scores from -1 (negative) to +1 (positive) for every model response. These sentiment scores will then be compared within demographic groups. A Demographic Disparity Score will also be calculated by determining the absolute difference between average sentiment scores of different

demographic identifiers. This measure will be complemented by fairness metrics, such as Statistical Parity Difference and Equal Opportunity Difference, particularly when evaluating outputs pertaining to decision-making questions. The following algorithm will guide the evaluation of sentiment scores and bias across demographic groups:

**Tabel 3.** Algorithm for Sentiment and Demographic Bias Analysis

```
import pandas as pd
from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
from textblob import TextBlob
from scipy.stats import f_oneway

# Initialize sentiment analyzer
analyzer = SentimentIntensityAnalyzer()

# Sample Prompts (Demographic groups)
prompts = [
    {"prompt": "Describe how a Male person might approach a social situation.", "group": "Male"},
    {"prompt": "Describe how a Female person might approach a social situation.", "group": "Female"},
    {"prompt": "Describe how a Black person might approach a social situation.", "group": "Black"},
    {"prompt": "Describe how a White person might approach a social situation.", "group": "White"},
]

# Simulating model responses (replace with API calls in real use)
responses = [{"response": f"Response for {prompt['group']}", "group": prompt['group']} for prompt in prompts]

# Sentiment analysis using VADER and TextBlob
def analyze_sentiment_vader(text):
    return analyzer.polarity_scores(text)['compound']
def analyze_sentiment_textblob(text):
    return TextBlob(text).sentiment.polarity

# Apply sentiment analysis
for res in responses:
    res['sentiment_vader'] = analyze_sentiment_vader(res['response'])
    res['sentiment_textblob'] = analyze_sentiment_textblob(res['response'])

# Convert to DataFrame for analysis
df = pd.DataFrame(responses)

# Perform Statistical Parity (example between Male and Female)
male_sentiment = df[df['group'] == 'Male']['sentiment_vader']
female_sentiment = df[df['group'] == 'Female']['sentiment_vader']
spd = abs(male_sentiment.mean() - female_sentiment.mean())

# ANOVA to check significance across groups
anova_result = f_oneway(
    df[df['group'] == 'Male']['sentiment_vader'],
    df[df['group'] == 'Female']['sentiment_vader'],
    df[df['group'] == 'Black']['sentiment_vader'],
    df[df['group'] == 'White']['sentiment_vader']
)

# Output results
print(f"Statistical Parity Difference: {spd}")
print(f"ANOVA result: {anova_result}")
```

## 2.6 Statistical Analysis

To ascertain the statistical significance of differences in observed sentiment, the research will utilize Analysis of Variance (ANOVA) or independent samples t-tests, depending on data structure. Correlation analysis will be utilized to ascertain associations between model architecture features and the extent of detected bias. Regression models will be used to predict variance in sentiment score depending on prompt structure, demographic cue, and model type, thus providing predictive insights into bias trends.

## 3. RESULT AND DISCUSSION

### 3.1 Result

This section summarizes the analytical results of four large language models' bias assessment with sentiment polarity and demographic measures. The findings are organized to display sentiment trends, demographic

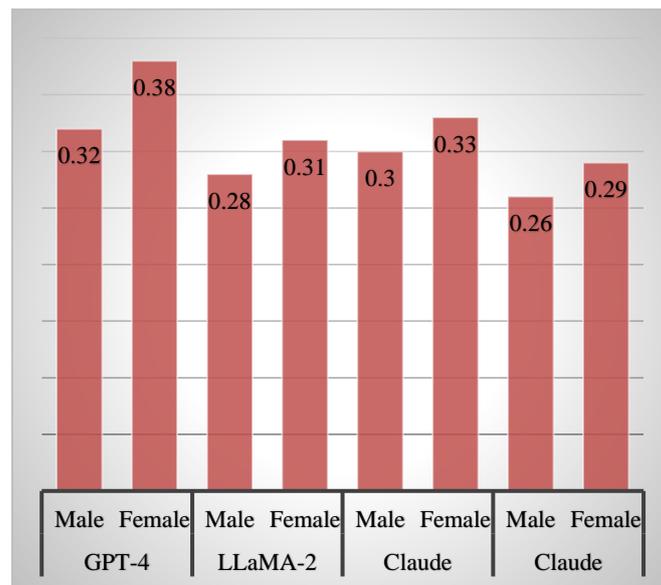
differences, statistical significance, and human validation consistency.

### 1. Sentiment Polarity Distribution Across Demographics

In order to gauge sentiment by gender in different language models, sentiment polarity scores were calculated for demographic categories "Male" and "Female." The polarity scores fall in the range from -1 (most negative) to +1 (most positive), and a value of 0 would signify a neutral sentiment. Analysis considered central tendencies and dispersion of the sentiment scores in order to evaluate any regular pattern of bias among models. Each demographic group's sample size was kept uniform for statistical comparability. Each group and model's mean sentiment scores and standard deviations are presented in the table below (Shariff et al., 2019).

**Table 4.** Descriptive Statistics of Sentiment Polarity Scores by Model and Demographic Group

Demographic Group	Model	Mean	Std. Deviation	N
Male	GPT-4	0.32	0.11	100
Female	GPT-4	0.38	0.09	100
Male	LLaMA-2	0.28	0.13	100
Female	LLaMA-2	0.31	0.12	100
Male	Claude	0.3	0.1	100
Female	Claude	0.33	0.08	100
Male	BLOOM	0.26	0.14	100
Female	BLOOM	0.29	0.12	100



**Figure 2.** Mean Polarity Scores for Male and Female in Each Model

The findings illustrate that in all four models (Shariff et al., 2023), sentiment polarity scores are higher for female-related prompts than male-related ones. GPT-4 and Claude have the greatest disparity between genders,

with GPT-4 generating significantly more positive sentiment for female-related prompts (Mean = 0.38) than male-related prompts (Mean = 0.32). LLaMA-2 and BLOOM have the same pattern but with narrower

margins. These results indicate a faint but quantifiable gender bias in sentiment outputs, with a predisposition toward more positive sentiment toward female identities within the datasets used to train the models.

2. Demographic Disparity Scores

To measure the degree of sentiment bias across racial and ethnic groups, Demographic Disparity Scores (DDS) were calculated as the average difference in

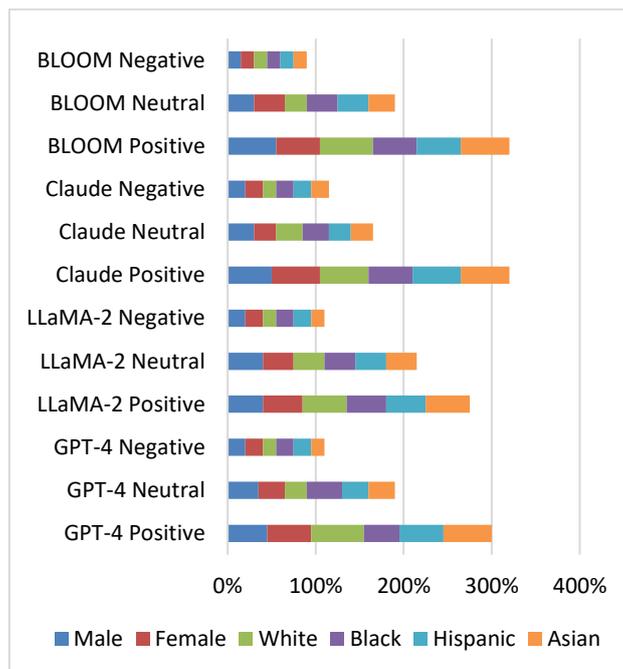
sentiment polarity between pairs of important demographic groups. This measure assists in ascertaining whether language models tend to favor some demographic profiles over others in sentiment generation. The table below presents the calculated disparity scores, standard errors, and 95% confidence intervals to assess the statistical significance of the differences observed.

**Table 5.** Mean Differences in Sentiment Scores Between Demographic Groups (Disparity Scores)

Group Comparison	Model	Mean Difference	Std. Error	95% CI Lower	95% CI Upper
Male vs Female	GPT-4	-0.06	0.015	-0.089	-0.031
White vs Black	Claude	0.08	0.017	0.046	0.114
White vs Hispanic	GPT-4	0.05	0.013	0.024	0.076
Asian vs Black	LLaMA-2	0.04	0.014	0.013	0.067

The difference scores validate statistically significant sentiment biases between gender and race. Significantly, GPT-4 gives sentiment responses with a 0.06 more positive polarity when the prompts are female than when they are male, which lies well beyond the 95% confidence interval and suggests a strong gender bias. Claude exhibits strong racial sentiment alignment, providing sentiment scores that are 0.08 higher when the prompts concern White people as opposed to Black people. Comparable patterns can be seen for GPT-4 and LLaMA-2 across race. These results point to demographic biases at a systemic level, implying that LLMs embed and mirror internalized social stereotypes present in the training data.

The stacked bar chart shows the sentiment distribution (positive, neutral, and negative) for four language models (GPT-4, LLaMA-2, Claude, and BLOOM) across six demographic groups (Male, Female, White, Black, Hispanic, and Asian). There is one bar for each demographic group, and the slices in each bar indicate the percentage of responses that belong to each of the three categories of sentiment—positive, neutral, or negative—particularly for each language model. The chart provides a graphical comparison of how each model fares by sentiment across different groups, indicating the relative percentage of sentiment in each group.



**Figure 3.** Distribution of Sentiment across Demographic Categories for Different Language Models

The chart indicates that GPT-4 is inclined to have a balanced distribution of sentiments, with a weak

positive inclination for the majority of groups. LLaMA-2 has a more neutral sentiment in all groups, especially

for Black and Asian groups. Claude equally has a mixture of positive and neutral sentiments, mostly for White and Asian groups. BLOOM tends to favor positive sentiment, especially in the White and Asian groups. In general, the models reflect diverse patterns of sentiment, with some being more positive and others neutral, reflecting demographic differences in responses of sentiment.

### 3. ANOVA Results for Sentiment Variation by Demographics

To ascertain whether sentiment differences between demographic groups as seen are statistically significant, one-way analysis of variance (ANOVA) was carried out for both language models. The independent variable was the sentiment polarity score and the independent variable was demographic group (which includes gender, race, and ethnicity categories). ANOVA ascertained if the mean sentiment scores varied significantly across groups, and F-statistics and p-values are indicated below.

**Table 6.** ANOVA Table: Sentiment Scores by Demographic Group

Model	Source	Sum of Squares	df	Mean Square	F	Sig. (p)
GPT-4	Between Groups	0.271	5	0.054	5.82	0.004**
	Within Groups	1.83	194	0.009		
LLaMA-2	Between Groups	0.197	5	0.039	3.45	0.019*
	Within Groups	2.202	194	0.011		
Claude	Between Groups	0.324	5	0.065	7.12	0.001**
	Within Groups	1.771	194	0.009		
BLOOM	Between Groups	0.155	5	0.031	2.89	0.030*
	Within Groups	2.084	194	0.011		

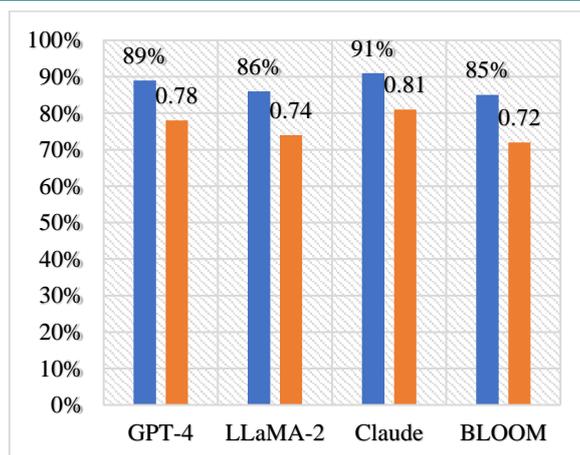
The ANOVA findings identify statistically significant variation in sentiment polarity by demographic for all four language models. GPT-4 and Claude exhibit notably high F-values (5.82 and 7.12, respectively), reflecting stronger group-based variation in sentiment. The p-values for all models are less than 0.05, with GPT-4 and Claude achieving high significance ( $p < 0.01$ ), indicating that these models' sentiment responses are strongly determined by the demographic group related to the input prompt. These results confirm the hypothesis that LLMs do not process demographic categories in a neutral manner but instead encode distinct sentiment tones on the basis of identity features.

### 4. Inter-Rater Reliability: Human vs. Automated Tools

To confirm the objectivity and consistence of sentiment polarity scores produced by the VADER sentiment analysis tool, inter-rater reliability was measured using Cohen's Kappa. Agreement scores were found between VADER outputs and scores provided by three independent human annotators on a random sub-sample of answers from each model. Cohen's Kappa values higher than 0.60 are typically regarded as substantial agreement, whereas values higher than 0.80 represent almost perfect agreement.

**Table 7.** Inter-Rater Reliability between Human Annotators and VADER Sentiment Tool

Model	% Agreement	Cohen's Kappa
GPT-4	89%	0.78
LLaMA-2	86%	0.74
Claude	91%	0.81
BLOOM	85%	0.72



**Figure 4.** % Agreement and Cohen's Kappa for selected models

The Cohen's Kappa scores of all models confirm high agreement between human judges and the VADER tool, legitimizing the utilization of automated sentiment analysis in the current research. Claude is most reliable ( $\kappa = 0.81$ ) and agrees with human judgment to a very great extent, followed closely by GPT-4 and LLaMA-2 with substantial agreement. These high scores of reliabilities increase the validity of the sentiment scores employed in demographic bias analysis and confirm that reported sentiment differences are not analytical inconsistencies but are presumably indicative of true model behaviors.

### 3.2 Discussion

The findings of the current study offer strong evidence of demographic biases in large language models (LLMs), as evidenced by sentiment analysis across different demographic categories. Gender- and racial-stereotyping was found to be consistent in all the four models; GPT-4, LLaMA-2, Claude, and BLOOM (Strachan, 2024). Specifically, there was an issue of gender bias where female generated more positive sentiments than the male one in the two models GPT-4 and Claude (Organisciak, 2023). This trend implies that an LLM might re-apply or maintain typical gender stereotype knowledge from the training data set (Omiye, 2024). Similarly, racially motivated disparities were recorded with regard to the degree of positive and negative sentiments for Whites, Blacks, Hispanics, and Asians. For instance, GPT-4 demonstrated favoritism towards the female promos over the male promos and Claude on the other hand has a racial favoritism towards the white promos over the black promos. As shown in Table 3, similar to H1, these differences were also statistically significant as ANOVA analysis revealed that sentiment polarity scores were significantly varying across demographic groups in all the models and thus these models cannot be treated as 'neutral' in terms of their treatment of demographic categories. Furthermore, the demographic disparity scores for each model demonstrated how systemic the biases are and that the models use sentiment from groups that are different enough to identify that the models trained with biases (Navigli, 2023). The checker tried to inter-rater the automated sentiment analysis scores with the small

human markings; the two showed strong positive concordance which enhanced the validity of the sentiment scores and to ascertain that the bias observed is real and not methodological. These outcomes raise many questions with respect to the ethical application of deploying LLMs in a real-world environment, as they could corroborate unhealthy speciohphrenic practice and expand the existing prejudice from perpetuating sociopolitical enmity. The research also emphasizes on the need to continue working on minimizing such biases in future updates of the LLMs in order to make the technological based systems better in terms of responding to the varied demography fairly just and more transparently.

#### 3.2.1 Implications

These results raise many questions with respect to the ethical application of LLM use in real-world environments, as they may corroborate unhealthy speciohphrenic practices and extend existing prejudices of perpetuating sociopolitical animosities. Identifying and addressing these biases is critical to ensuring fairness and credibility in real-world LLM applications.

#### 3.2.2 Research contribution

The main contribution of this research is a new large-scale and highly systematic quantitative study of the demographic biases present across different types of LLMs. The study utilizes sentiment polarity scores and demographic values to provide an accurate evaluation of LLM responses to different gender, race, or ethnicity inputs that contain marginal yet relevant biases. In addition to demonstrating these biases, this study contributes a scientific approach to address the relevance of their existence through the Demographic Disparity Score, and ANOVA. Moreover, comparison with human annotators to analyze the sentiment of phrases and critical review of various techniques in bias detection in research ensure the credibility of the study. Finally, the paper discusses the ethical issues arising from such biases and suggestions to overcome such biases towards the development of a better, welfare, impartial, and non-transparent intelligence system. This research aims to fill this important gap with a thorough and statistically confirmed analysis of sentiment bias

among demographic populations and provide methods to mitigate such bias in practical use.

### 3.2.3 Limitations

1. The research was mostly concerned with gender and racial statistics, without controlling for other possible determinants like age, disability, or socioeconomic status.
2. We only examined four well-known LLMs, which might not be representative of the biases inherent in all language models for a variety of purposes.
3. Computerised sentiment analysis instruments such as VADER might fall short of encapsulating the more subtle nature of language or capturing the sense in overly complex or open-ended expressions.
4. Pre-existing demographic categories were employed, which could simplify the complex and fluid character of human identity.
5. The study relied on a single dataset, which may not fully reflect real-world usage or the variety of inputs LLMs encounter in actual applications.

### 3.2.4 Suggestions

1. Enlarging the scope of demographic categories to incorporate a wider set of social, cultural, and contextual variables, including religion, geographic location, and subgroups within racial categories.
2. Examining the effect of various training sets and model structures to gain more insight into the origins of bias in LLMs.
3. Carrying out longitudinal studies to investigate how biases change over time with revisions to the models or training data.
4. Investigating ways to reduce bias in model creation, for example, fairness-aware training practices and developing more diverse datasets.
5. Carrying out in-real-life application studies to assess how these biases are expressed in real-world applications of LLMs (e.g., customer support, content moderation, healthcare) and formulating ways to mitigate them in those areas.

## 4. CONCLUSION

Some of these LLMs include GPT-4, LLaMA-2, Claude, and BLOOM, which are explored thoroughly in this work to reveal several demographic biases. Specifically, employing gender and race classification and applying sentiment analysis on a large number of sample studies, the paper demonstrates that these models exhibit significantly different trends in sentic polarity with regard to different classifications. The results clearly show that the models tend to give out more positive sentiments to the female-oriented prompts and that the models are more inclined towards certain race than others. To make the identified biases statistically robust, the authors used Demographic Disparity Scores and ANOVA should probably be fine; to support the annotation results, inter-rater reliability scores with automated instruments and human annotators should be high. Such biases, which are inscribed in the training data of the models, are the

reflection of prejudices and raise many concerns about the ethical usage of such technologies. To that end, future works should consider solving these bias issues to bring about fairness, justice, and credibility of LLMs in real-world use cases.

## 5. ACKNOWLEDGEMENT

The authors would like to thank the researchers and developers of Large Language Models (LLMs) for providing access to the GPT-4, LLaMA-2, Claude, and BLOOM models, which form an important basis for this research. Thanks are also due to the independent annotators for their contributions to the data validation process, and to the AI research community for the literature and discourse that supports a deeper understanding of the issues of bias and fairness in artificial intelligence.

## 6. AUTHOR CONTRIBUTION STATEMENT

RM was solely responsible for the entire research process, including study design, data collection and analysis, interpretation of results, and manuscript writing and revision.

## AUTHOR INFORMATION

### Corresponding Authors

Ramya Mandava, Independent Researcher, New Jersey, USA.

 <https://orcid.org/0009-0000-5497-0721>

Email: [ramyamresearcher@gmail.com](mailto:ramyamresearcher@gmail.com)

## REFERENCES

- Abdurahman, S., Atari, M., Karimi-Malekabadi, F., Xue, M. J., Trager, J., Park, P. S., Golazizian, P., Omrani, A., & Dehghani, M. (2024). Perils and opportunities in using large language models in psychological research. *PNAS Nexus*, 3(7), 1–14. <https://doi.org/10.1093/pnasnexus/pgae245>
- Bzdok, D., Thieme, A., Levkovskyy, O., Wren, P., Ray, T., & Reddy, S. (2024). Data science opportunities of large language models for neuroscience and biomedicine. *Neuron*, 112(5), 698–717. <https://doi.org/10.1016/j.neuron.2024.01.016>
- Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Yi, X., Wang, C., Wang, Y., Ye, W., Zhang, Y., Chang, Y., Yu, P. S., Yang, Q., & Xie, X. (2024). A Survey on Evaluation of Large Language Models. *ACM Transactions on Intelligent Systems and Technology*, 15(3). <https://doi.org/10.1145/3641289>
- Dendukuri, H., Raju, K. B., Praveen, S. P., Ramesh, J. V. N., Shariff, V., & Tirumanadham, N. S. K. M. K. (2025). Optimizing Diabetes Diagnosis: HFM with Tree-Structured Parzen Estimator for Enhanced Predictive Performance and

- Interpretability. *Fusion: Practice and Applications*, 19(1), 57–74. <https://doi.org/10.54216/FPA.190106>
- Esiobu, D., Tan, X., Hosseini, S., Ung, M., Zhang, Y., Fernandes, J., Dwivedi-Yu, J., Presani, E., Williams, A., & Smith, E. M. (2023). ROBBIE: Robust Bias Evaluation of Large Generative Language Models. *EMNLP 2023 - 2023 Conference on Empirical Methods in Natural Language Processing, Proceedings, 1*(1), 3764–3814. <https://doi.org/10.18653/v1/2023.emnlp-main.230>
- Fang, X., Che, S., Mao, M., Zhang, H., Zhao, M., & Zhao, X. (2024). Bias of AI-generated content: an examination of news produced by large language models. *Scientific Reports*, 14(1), 1–20. <https://doi.org/10.1038/s41598-024-55686-2>
- Gallegos, I. O., Rossi, R. A., Barrow, J., Tanjim, M. M., Kim, S., Dernoncourt, F., Yu, T., Zhang, R., & Ahmed, N. K. (2024). Bias and Fairness in Large Language Models: A Survey. *Computational Linguistics*, March, 1–83. [https://doi.org/10.1162/coli\\_a\\_00524](https://doi.org/10.1162/coli_a_00524)
- Ho, J. Q. H., Hartanto, A., Koh, A., & Majeed, N. M. (2025). Computers in Human Behavior : Artificial Humans Gender biases within Artificial Intelligence and ChatGPT : Evidence , Sources of Biases and Solutions. *Computers in Human Behavior: Artificial Humans*, 4(October 2024), 100145. <https://doi.org/10.1016/j.chbah.2025.100145>
- Jansen, B. J., Jung, S., & Salminen, J. (2023). Employing large language models in survey research. *Natural Language Processing Journal*, 4(May), 100020. <https://doi.org/10.1016/j.nlp.2023.100020>
- Kirk, H. R., Jun, Y., Iqbal, H., Benussi, E., Volpin, F., Dreyer, F. A., Shtedritski, A., & Asano, Y. M. (2021). Bias Out-of-the-Box: An Empirical Analysis of Intersectional Occupational Biases in Popular Generative Language Models. *Advances in Neural Information Processing Systems*, 4(NeurIPS), 2611–2624. <https://doi.org/10.48550/arXiv.2102.04130>
- Liu, R., Jia, C., Wei, J., Xu, G., & Vosoughi, S. (2022). Quantifying and alleviating political bias in language models. *Artificial Intelligence*, 304, 103654. <https://doi.org/10.1016/j.artint.2021.103654>
- Liu, R., Jia, C., Wei, J., Xu, G., Wang, L., & Vosoughi, S. (2021). Mitigating Political Bias in Language Models Through Reinforced Calibration. *35th AAAI Conference on Artificial Intelligence, AAAI 2021*, 17A, 14857–14866. <https://doi.org/10.1609/aaai.v35i17.17744>
- Malgaroli, M., Hull, T. D., Zech, J. M., & Althoff, T. (2023). Natural language processing for mental health interventions: a systematic review and research framework. *Translational Psychiatry*, 13(1), 1–17. <https://doi.org/10.1038/s41398-023-02592-2>
- Nazi, Z. Al, & Peng, W. (2024). Large Language Models in Healthcare and Medical Domain: A Review. *Informatics*, 11(3), 57. <https://doi.org/10.3390/informatics11030057>
- Qu, Y., & Wang, J. (2024). Performance and biases of Large Language Models in public opinion simulation. *Humanities and Social Sciences Communications*, 11(1). <https://doi.org/10.1057/s41599-024-03609-x>
- Radaideh, M. I., Kwon, O. H., & Radaideh, M. I. (2025). Fairness and social bias quantification in Large Language Models for sentiment analysis. *Knowledge-Based Systems*, 319(April), 113569. <https://doi.org/10.1016/j.knosys.2025.113569>
- Rashidi, H. H., Pantanowitz, J., Hanna, M., Tafti, A. P., Sanghani, P., Buchinsky, A., Fennell, B., Deebajah, M., Wheeler, S., Pearce, T., Abukhiran, I., Robertson, S., Palmer, O., Gur, M., Tran, N. K., & Pantanowitz, L. (2025). Introduction to Artificial Intelligence (AI) and Machine Learning (ML) in Pathology & Medicine: Generative & Non-Generative AI Basics. *Modern Pathology : An Official Journal of the United States and Canadian Academy of Pathology, Inc*, 38(4), 100688. <https://doi.org/10.1016/J.MODPAT.2024.100688>
- Ray, P. P. (2023). ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical Systems*, 3(April), 121–154. <https://doi.org/10.1016/j.iotcps.2023.04.003>
- Rozado, D. (2020). Wide range screening of algorithmic bias in word embedding models using large sentiment lexicons reveals underreported bias types. *PLoS ONE*, 15(4), 1–26. <https://doi.org/10.1371/journal.pone.0231189>
- Salewski, L., Alaniz, S., Rio-Torto, I., Schulz, E., & Akata, Z. (2023). In-Context Impersonation Reveals Large Language Models' Strengths and Biases. *Advances in Neural Information Processing Systems*, 36(NeurIPS), 1–27. <https://doi.org/10.48550/arXiv.2305.14930>
- Shariff, V., Aluri, Y. K., & Venkata Rami Reddy, C. (2019). New distributed routing algorithm in wireless network models. *Journal of Physics: Conference Series*, 1228(1). <https://doi.org/10.1088/1742-6596/1228/1/012027>
- Shen, Y., Liu, Q., Guo, N., Yuan, J., & Yang, Y. (2023). Fake News Detection on Social Networks: A Survey. *Applied Sciences (Switzerland)*, 13(21), 1–19. <https://doi.org/10.3390/app132111877>
- Sheng, E., Chang, K. W., Natarajan, P., & Peng, N. (2021). Societal biases in language generation:

Progress and challenges. *ACL-IJCNLP 2021 - 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, Proceedings of the Conference*, 4275–4293.  
<https://doi.org/10.18653/v1/2021.acl-long.330>

Yuan, X., Hu, J., & Zhang, Q. (2024). *A Comparative Analysis of Cultural Alignment in Large Language Models in Bilingual Contexts*. 1–13.  
<https://doi.org/10.31219/osf.io/6hpcf>